# Meta Reinforcement Learning

Kate Rakelly 11/13/19

#### Questions we seek to answer

Motivation: What problem is meta-RL trying to solve?

**Context**: What is the connection to other problems in RL?

**Solutions**: What are solution methods for meta-RL and their limitations?

**Open Problems**: What are the open problems in meta-RL?

# Meta-learning problem statement

#### supervised learning



reinforcement learning





Robot art by Matt Spangler, mattspangler.com

### Meta-RL problem statement

**Regular RL**: learn policy for single task

$$\theta^{\star} = \arg \max_{\theta} E_{\pi_{\theta}(\tau)}[R(\tau)]$$
$$= f_{\mathrm{RL}}(\mathcal{M})$$
$$\bigwedge_{\mathrm{MDP}}$$



Meta-RL: learn adaptation rule





 $\mathcal{M}_1$ 

 $\mathcal{M}_{test}$ 

### Relation to goal-conditioned policies



Meta-RL can be viewed as a goal-conditioned policy where the task information is inferred from *experience* 

Task information could be about the dynamics or reward functions

Rewards are a strict generalization of goals

### Relation to goal-conditioned policies



Q: What is an example of a reward function that can't be expressed as a goal state?

A: E.g., seek while avoiding, action penalties

### Adaptation

$$\theta^{\star} = \arg \max_{\theta} \sum_{i=1}^{n} E_{\pi_{\phi_i}(\tau)}[R(\tau)]$$

where 
$$\phi_i = f_{\theta}(\mathcal{M}_i)$$



#### What should the adaptation procedure do?

- **Explore**: Collect the most informative data
- Adapt: Use that data to obtain the optimal policy

### General meta-RL algorithm outline



Different algorithms:

- Choice of function f
- Choice of loss function L

# **Solution Methods**



Persist the hidden state across episode boundaries for continued adaptation!

Duan et al. 2016, Wang et al. 2016. Heess et al. 2015. Fig adapted from Duan et al. 2016

#### Solution #1: recurrence

while training: for *i* in tasks: initialize hidden state  $h_0 = 0$ for *t* in timesteps:

1. sample 1 transition  $\mathcal{D}_i = \mathcal{D}_i \cup \{(s_t, a_t, s_{t+1}, r_t)\}$  from  $\pi_{h_t}$ 

2. update policy hidden state  $\mathbf{h_{t+1}} = f_{\theta}(\mathbf{h_t}, s_t, a_t, s_{t+1}, r_t)$ 

update policy parameters  $\theta \leftarrow \theta - \nabla_{\theta} \sum_{i} \mathcal{L}_{i}(\mathcal{D}_{i}, \pi_{\mathbf{h}})$ 

#### Solution #1: recurrence

#### Pro: general, expressive

There exists an RNN that can compute any function

#### Con: not consistent

What does it mean for adaptation to be "consistent"?

Will converge to the optimal policy given enough data



### Solution #1: recurrence











(a) Labryinth I-maze

(b) Illustrative Episode

Duan et al 2016, Wang et al. 2016

### Wait, what if we just fine-tune?



is pretraining a *type* of meta-learning? better features = faster learning of new task!

Sample inefficient, prone to overfitting, and is particularly difficult in RL

Slide adapted from Sergey Levine

 $abla \mathcal{L}_3$ 

 $abla \mathcal{L}_2$ 

Learn a parameter initialization from which fine-tuning for a new task works!

 $abla \mathcal{L}$ 

 $\theta_1^*$ 

---- meta-learning ---- learning/adaptation

 $\theta^*$ 

· 02



 $\theta^{\star}$ 

=



n

where  $\phi_i = f_{\theta}(\mathcal{M}_i)$ 

PG

argmax

 $E_{\pi_{\phi_i}(\tau)}[R(\tau)]$ 

while training: for *i* in tasks:

1. sample k episodes  $\mathcal{D}_i = \{(s, a, s', r)\}_{1:k}$  from  $\pi_{\theta}$ 

2. compute adapted parameters  $\phi_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(\pi_{\theta}, \mathcal{D}_i)$ 

3. sample k episodes  $\mathcal{D}'_i = \{(s, a, s', r)_{1:k}\}$  from  $\pi_{\phi}$ 

update policy parameters  $\theta \leftarrow \theta - \nabla_{\theta} \sum_{i} \mathcal{L}_{i}(\mathcal{D}'_{i}, \pi_{\phi_{i}})$ 

Requires second order derivatives!

How exploration is learned automatically



Pre-update parameters receive credit for producing good exploration trajectories T $\pi_{\theta}$ U $\pi_{\theta'}$  $\pi_{\theta'}$  $\tau$ 

Causal relationship between pre and post-update trajectories is taken into account

R'

View this as a "return" that encourages gradient alignment

**Pro: consistent!** 

Con: not as expressive



#### Q: When could the optimization strategy be less expressive than the recurrent strategy?

Example: when no rewards are collected, adaptation will not change the policy, even though this data gives information about which states to avoid



Exploring in a sparse reward setting



Fig adapted from Rothfuss et al. 2018

Cheetah running forward and back after 1 gradient step



Fig adapted from Finn et al. 2017

# Meta-RL on robotic systems

# Meta-imitation learning

#### Demonstration



#### 1-shot imitation



Figure adapted from BAIR Blog Post: One-Shot Imitation from Watching Videos

# Meta-imitation learning

Test: perform task given single **robot demo** Training: run **behavior cloning** for adaptation



learn how to infer a policy \_\_\_\_\_\_\_ from one demonstration

Test time provide 1 demo with new object





$$\phi_i = heta - lpha 
abla_ heta \sum_t ||\pi_ heta(o_t) - a_t^*||^2$$

## Meta-imitation learning from human demos

demonstration



1-shot imitation



Figure adapted from BAIR Blog Post: One-Shot Imitation from Watching Videos

# Meta-imitation learning from humans

Test: perform task given single **human demo** Training: **learn a loss function** that adapts policy



where  $\phi_i = \underbrace{f_{\theta}}_{i} \mathcal{M}_i$ 

n

arg max

PG

Learned loss

 $E_{\pi_{\phi_i}(\tau)}[R(\tau)]$ 

$$\phi = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\psi}(\theta, \mathbf{d}^h)$$



Supervised by **paired robot-human demos** only during meta-training!

### Model-Based meta-RL

1. run base policy  $\pi_0(\mathbf{a}_t|\mathbf{s}_t)$  (e.g., random policy) to collect  $\mathcal{D} = \{(\mathbf{s}, \mathbf{a}, \mathbf{s}')_i\}$ 

- 2. learn dynamics model  $f(\mathbf{s}, \mathbf{a})$  to minimize  $\sum_i ||f(\mathbf{s}_i, \mathbf{a}_i) \mathbf{s}'_i||^2$
- 3. plan through  $f(\mathbf{s}, \mathbf{a})$  to choose actions



What if the system dynamics change?

- Low battery
- Malfunction
- Different terrain

Re-train model? :(



Figure adapted from Anusha Nagabandi

#### Model-Based meta-RL





### Aside: POMDPs

# Example: incomplete sensor data



"That Way We Go" by Matt Spangler

state is unobserved (hidden)  $h_t$  $n_{t+2}$  $n_{t+1}$  $o_t$  $o_{t+1}$  $o_{t+2}$ observation gives  $a_{t+1}$  $a_t$ incomplete information about the state

### The POMDP view of meta-RL



...as a POMDP

$$h_t = (s_t, task)$$
  $o_t = (s_t, r_t)$ 



Two approaches to solve: 1) policy with memory (RNN)

2) explicit state estimation

### Model belief over latent task variables

POMDP for unobserved state

#### **POMDP** for unobserved task



a = "left", s = S0, r = 0

a = "left", s = S0, r = 0

### Model belief over latent task variables

**POMDP** for unobserved state

#### **POMDP** for unobserved task



a ="left", s = S0, r = 0

a = "left", s = S0, r = 0



### Solution #3: posterior sampling in action





See Control as Inference (Levine 2018) for justification of thinking of Q as a pseudo-likelihood



#### Aside: Soft Actor-Critic (SAC)

"Soft": Maximize rewards \*and\* entropy of the policy (higher entropy policies explore better)

$$J(\pi) = \sum_{t=0}^{T} \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_{\pi}} \left[ r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi(\cdot | \mathbf{s}_t)) \right]$$

"Actor-Critic": Model \*both\* the actor (aka the policy) and the critic (aka the Q-function)

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right]$$
$$J_\pi(\phi) = \mathbb{E}_{s_t, a_t} [Q_\theta(s_t, a_t) + \alpha \mathcal{H}(\pi_\phi(\cdot | s_t))]$$

Much more **sample efficient** than on-policy algs.



Dclaw robot turns valve from pixels

SAC Haarnoja et al. 2018, Control as Inference Tutorial. Levine 2018, SAC BAIR Blog Post 2019

#### Soft Actor-Critic





### Meta-RL experimental domains



Simulated via MuJoCo (Todorov et al. 2012), tasks proposed by (Finn et al. 2017, Rothfuss et al. 2019)



ProMP (Rothfuss et al. 2019), MAML (Finn et al. 2017), RL2 (Duan et al. 2016)



ProMP (Rothfuss et al. 2019), MAML (Finn et al. 2017), RL2 (Duan et al. 2016)

### two views of meta-RL

#### Mechanistic view

- Deep neural network model that can read in an entire dataset and make predictions for new datapoints
- Training this network uses a meta-dataset, which itself consists of many datasets, each for a different task

#### Probabilistic view

- Extract prior information from a set of (metatraining) tasks that allows efficient learning of new tasks
- Learning a new task uses this prior and (small) training set to infer most likely posterior parameters

### Summary



Slide adapted from Sergey Levine and Chelsea Finn

# Frontiers

### Where do tasks come from?

Idea: generate self-supervised tasks and use them during meta-training



Skills should be high entropy

Point robot learns to explore different areas after the hallway

states



Ant learns to run in different directions, jump, and flip



Limitations
Assumption that skills shouldn't depend on action not always valid
Distribution shift meta-train -> meta-test

# How to explore efficiently in a new task?

Learn exploration strategies better...









#### Bias exploration with extra information...



human -provided demo



Robot attempt #1, w/ only demo info



Robot attempt #2, w/ demo + reward info

Gupta et al. 2018, Rakelly et al. 2019, Zhou et al. 2019

### Online meta-learning

Meta-training tasks are presented in a sequence rather than a batch



#### Summary

Meta-RL finds an adaptation procedure that can quickly adapt the policy to a new task

Three main solution classes: RNN, optimization, task-belief and several learning paradigms: model-free (on and off policy), model-based, imitation learning

Connection to goal-conditioned RL and POMDPs

Some open problems (there are more!): better exploration, defining task distributions, meta-learning online

#### References

#### **Recurrent meta-RL**

Learning to Reinforcement Learn, Wang et al. 2016 Fast Reinforcement Learning by Slow Reinforcement Learning, Duan et al. 2016 Memory-Based Control with Recurrent Neural Networks, Heess et al. 2015

#### **Optimization-based meta-RL**

Model-Agnostic Meta-Learning, Finn et al. 2017 Proximal Meta-Policy Search, Rothfuss et al. 2018

#### Optimization-based meta-RL + imitation learning

One-Shot Visual Imitation Learning via Meta-Learning, Yu et al. 2017 One-Shot Imitation from Observing Humans via Domain-Adaptive Meta-Learning, Yu et al. 2018

#### Model-based meta-RL

Learning to Adapt in Dynamic, Real-World Environments through Meta-Reinforcement Learning, Nagabandi et al. 2019

#### Off-policy meta-RL

Soft Actor-Critic, Haarnoja et al. 2018 Control as Inference, Levine 2018. Efficient Off-Policy Meta-RL via Probabilistic Context Variables, Rakelly et al. 2019

#### References

#### **Open Problems**

Diversity is All You Need: Learning Skills without a Reward Function, Eysenbach et al. 2018 Unsupervised Meta-learning for RL, Gupta et al. 2018 Meta-Reinforcement Learning of Structured Exploration Strategies, Gupta et al. 2018 Watch, Try, Learn, Meta-Learning from Demonstrations and Reward, Zhou et al. 2019 Online Meta-Learning, Finn et al. 2019

#### Slides and Figures

Some slides adapted from Meta-Learning Tutorial at ICML 2019, Finn and Levine Robot illustrations by Matt Spangler, mattspangler.com