Meta-Learning

CS 294-112: Deep Reinforcement Learning

Sergey Levine

Class Notes

- 1. Two weeks until the project milestone!
- 2. Guest lectures start next week, be sure to attend!
- 3. Today: part 1: meta-learning
- 4. Today: part 2: parallelism

How can we frame transfer learning problems? No single solution! Survey of various recent research papers

- 1. "Forward" transfer: train on one task, transfer to a new task
 - a) Just try it and hope for the best
 - b) Finetune on the new task
 - c) Architectures for transfer: progressive networks
 - d) Randomize source task domain
- 2. Multi-task transfer: train on many tasks, transfer to a new task
 - a) Model-based reinforcement learning
 - b) Model distillation
 - c) Contextual policies
 - d) Modular policy networks
- 3. Multi-task meta-learning: learn to learn from many tasks
 - a) RNN-based meta-learning
 - b) Gradient-based meta-learning

So far...

- Forward transfer: source domain to target domain
 - Diversity is good! The more varied the training, the more likely transfer is to succeed
- Multi-task learning: even more variety
 - No longer training on the same kind of task
 - But more variety = more likely to succeed at transfer
- How do we represent transfer knowledge?
 - Model (as in model-based RL): rules of physics are conserved across tasks
 - Policies requires finetuning, but closer to what we want to accomplish
 - What about *learning methods*?

What is meta-learning?

- If you've learned 100 tasks already, can you figure out how to *learn* more efficiently?
 - Now having multiple tasks is a huge advantage!
- Meta-learning = *learning to learn*
- In practice, very closely related to multi-task learning
- Many formulations
 - Learning an optimizer
 - Learning an RNN that ingests experience
 - Learning a representation



Why is meta-learning a good idea?

- Deep reinforcement learning, especially model-free, requires a huge number of samples
- If we can meta-learn a faster reinforcement learner, we can learn new tasks efficiently!
- What can a *meta-learned* learner do differently?
 - Explore more intelligently
 - Avoid trying actions that are know to be useless
 - Acquire the right features more quickly

Meta-learning with supervised learning

training data

test set





image credit: Ravi & Larochelle '17

Meta-learning with supervised learning





supervised learning:
$$f(x) \rightarrow y$$

 $f \qquad \uparrow$
input (e.g., image) output (e.g., label)

supervised meta-learning: $f(\mathcal{D}_{\text{train}}, x) \to y$ ftraining set

- How to read in training set?
 - Many options, RNNs can work
 - More on this later

The meta-learning problem in RL

supervised meta-learning: $f(\mathcal{D}_{\text{train}}, x) \to y$

reinforcement meta-learning (for example...): $f(\mathcal{D}_{\text{train}}, s) \to a$ $f f \uparrow \uparrow$ (recent experience state output (e.g., action)

$$\mathcal{D}_{\text{train}} = \{s_1, a_1, r_1, \dots, a_N, s_N, r_N\}$$



Meta-learning in RL with memory

"water maze" task



Heess et al., "Memory-based control with recurrent neural networks."







Duan et al., "RL2: Fast Reinforcement Learning via Slow Reinforcement Learning"

Connection to contextual policies

contextual policy: $\pi_{\theta}(\mathbf{a}|\mathbf{s},\omega)$





 $\omega :$ stack location

 $\omega :$ walking direction

just contextual policies, with experience as context

Back to representations...



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



Finn et al., "Model-Agnostic Meta-Learning"

What did we just do??

supervised learning: $f(x) \to y$

supervised meta-learning: $f(\mathcal{D}_{\text{train}}, x) \to y$

model-agnostic meta-learning: $f_{\text{MAML}}(\mathcal{D}_{\text{train}}, x) \to y$

$$f_{\text{MAML}}(\mathcal{D}_{\text{train}}, x) = f_{\theta'}(x)$$

$$\theta' = \theta - \alpha \sum_{(x,y) \in \mathcal{D}_{\text{train}}} \nabla_{\theta} \mathcal{L}(f_{\theta}(x), y)$$

Just another computation graph... Can implement with any autodiff package (e.g., TensorFlow) But has favorable inductive bias... Model-agnostic meta-learning: accelerating PG

after MAML training





after 1 gradient step

(forward reward)



after 1 gradient step

(backward reward)



 $\begin{array}{c} & & & \text{meta-learning} \\ & & & \text{learning/adaptation} \\ & & & \nabla \mathcal{L}_3 \\ & & & \nabla \mathcal{L}_2 \\ & & & \nabla \mathcal{L}_1 \\ & & & \nabla \mathcal{L}_2 \\ & & & & \partial_3^* \\ & & & & \partial_1^* \\ \end{array}$

Model-agnostic meta-learning: accelerating PG

after MAML training





after 1 gradient step

(backward reward)



 $\begin{array}{c} & ---- \text{meta-learning} \\ & ---- \text{learning/adaptation} \\ & \nabla \mathcal{L}_3 \\ & \nabla \mathcal{L}_2 \\ & \nabla \mathcal{L}_1 \\ & & \nabla \mathcal{L}_2 \\ & & & \Theta_3^* \\ & & & & & \Theta_2^* \end{array}$

after 1 gradient step

(forward reward)





Meta-learning summary & open problems

- Meta-learning = learning to learn
- Supervised meta-learning = supervised learning with datapoints that are entire datasets
- RL meta-learning with RNN policies
 - Ingest past experience with RNN
 - Simply run forward pass at test time to "learn"
 - Just contextual policies (no actual learning)
- Model-agnostic meta-learning
 - Use gradient descent (e.g., policy gradient) learning rule
 - Conceptually not that different
 - ...but can accelerate standard RL algorithms (e.g., learn in one iteration of PG)

Meta-learning summary & open problems

- The promise of meta-learning: use past experience to simply acquire a much more efficient deep RL algorithm
- The reality of meta-learning: mostly works well on smaller problems
- ...but getting better all the time
- Main limitations
 - RNN policies are extremely hard to train, and likely not scalable
 - Model-agnostic meta-learning presents a tough optimization problem
 - Designing the right task distribution is hard
 - Generally very sensitive to task distribution (meta-overfitting)

Parallelism in RL

Overview

- 1. We learned about a number of policy search methods
- 2. These algorithms have all been *sequential*
- 3. Is there a natural way to parallelize RL algorithms?
 - Experience sampling vs learning
 - Multiple learning threads
 - Multiple experience collection threads

Today's Lecture

- 1. What can we parallelize?
- 2. Case studies: specific parallel RL methods
- 3. Tradeoffs & considerations
- Goals
 - Understand the high-level anatomy of reinforcement learning algorithms
 - Understand standard strategies for parallelization
 - Tradeoffs of different parallel methods

High-level RL schematic



Which parts are slow?



Which parts can we parallelize?



Helps to group data generation and training (worker generates data, computes gradients, and gradients are pooled)

High-level decisions

- 1. Online or batch-mode?
- 2. Synchronous or asynchronous?





Relationship to parallelized SGD





- 1. Parallelizing model/critic/actor training typically involves parallelizing SGD
- 2. Simple parallel SGD:
 - 1. Each worker has a different slice of data
 - 2. Each worker computes gradients, sums them, sends to parameter server
 - 3. Parameter server sums gradients from all workers and sends back new parameters
- 3. Mathematically equivalent to SGD, but not asynchronous (communication delays)
- 4. Async SGD typically does not achieve perfect parallelism, but lack of locks can make it much faster
- 5. Somewhat problem dependent

Simple example: sample parallelism with PG

- 1. collect samples $\tau_i = \{\mathbf{s}_1^i, \mathbf{a}_1^i, \dots, \mathbf{s}_T^i, \mathbf{a}_T^i\}$ by running $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ N times
- 2. compute $r_i = r(\tau_i)$
- 3. compute $\nabla_i = \left(\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i)\right) (r_i b)$
- 4. update: $\theta \leftarrow \theta + \alpha \sum_i \nabla_i$



Simple example: sample parallelism with PG

- 1. collect samples $\tau_i = \{\mathbf{s}_1^i, \mathbf{a}_1^i, \dots, \mathbf{s}_T^i, \mathbf{a}_T^i\}$ by running $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ N times
- 2. compute $r_i = r(\tau_i)$
- 3. compute $\nabla_i = \left(\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i)\right) (r_i b)$
- 4. update: $\theta \leftarrow \theta + \alpha \sum_i \nabla_i$



Simple example: sample parallelism with PG

- 1. collect samples $\tau_i = \{\mathbf{s}_1^i, \mathbf{a}_1^i, \dots, \mathbf{s}_T^i, \mathbf{a}_T^i\}$ by running $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ N times
- 2. compute $r_i = r(\tau_i)$
- 3. compute $\nabla_i = \left(\sum_t \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i | \mathbf{s}_t^i)\right) (r_i b)$
- 4. update: $\theta \leftarrow \theta + \alpha \sum_i \nabla_i$





What if we add a critic?

- 1. collect samples $\tau_i = \{\mathbf{s}_1^i, \mathbf{a}_1^i, \dots, \mathbf{s}_T^i, \mathbf{a}_T^i\}$ by running $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ N times
- 2. compute $r_i = r(\tau_i)$

3. update $\hat{A}_{\phi}(\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i})$ with regression to target values \leftarrow see John's actor-critic lecture 4. compute $\nabla_{i} = \left(\sum_{t} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t}^{i} | \mathbf{s}_{t}^{i})\right) \hat{A}_{\phi}(\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i})$

5. update: $\theta \leftarrow \theta + \alpha \sum_i \nabla_i$



What if we add a critic?

- 1. collect samples $\tau_i = {\mathbf{s}_1^i, \mathbf{a}_1^i, \dots, \mathbf{s}_T^i, \mathbf{a}_T^i}$ by running $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ N times
- 2. compute $r_i = r(\tau_i)$

3. update $\hat{A}_{\phi}(\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i})$ with regression to target values \leftarrow see John's actor-critic lecture 4. compute $\nabla_{i} = \left(\sum_{t} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t}^{i} | \mathbf{s}_{t}^{i})\right) \hat{A}_{\phi}(\mathbf{s}_{t}^{i}, \mathbf{a}_{t}^{i})$

5. update: $\theta \leftarrow \theta + \alpha \sum_i \nabla_i$



What if we run online?

- 1. collect sample $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i)$ by running $\pi_{\theta}(\mathbf{a}|\mathbf{s})$ for 1 step
- 2. compute $r_i = r(\mathbf{s}_i, \mathbf{a}_i)$
- 3. update $\hat{A}_{\phi}(\mathbf{s}_t^i, \mathbf{a}_t^i)$ with regression to target values

4. compute
$$\nabla_i = \nabla_\theta \log \pi_\theta(\mathbf{a}^i | \mathbf{s}^i) \hat{A}_\phi(\mathbf{s}^i, \mathbf{a}^i)$$

5. update: $\theta \leftarrow \theta + \alpha \sum_i \nabla_i$

only the parameter update requires synchronization (actor + critic params)



Actor-critic algorithm: A3C

- 1. collect sample $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i)$ by running $\pi_{\theta}(\mathbf{a}|\mathbf{s})$ for 1 step
- 2. compute $r_i = r(\mathbf{s}_i, \mathbf{a}_i)$
- 3. update $\hat{A}_{\phi}(\mathbf{s}_t^i, \mathbf{a}_t^i)$ with regression to target values
- 4. compute $\nabla_i = \nabla_\theta \log \pi_\theta(\mathbf{a}^i | \mathbf{s}^i) \hat{A}_\phi(\mathbf{s}^i, \mathbf{a}^i)$
- 5. update: $\theta \leftarrow \theta + \alpha \sum_i \nabla_i$ (only do this every *n* steps)



- Some differences vs DQN, DDPG, etc:
 - No replay buffer, instead rely on diversity of samples from different workers to decorrelate
 - Some variability in exploration between workers
- Pro: generally much faster in terms of wall clock
- Con: generally must slower in terms of # of samples (more on this later...)

Mnih et al. '16

Actor-critic algorithm: A3C



Model-based algorithms: parallel GPS

- 1. get N samples τ_i by running $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ N times for each initial state \mathbf{s}_0^j
- 2. fit local models for each initial state
- 3. use LQR to get updated local policies $p_j(\mathbf{a}_t|\mathbf{s}_t)$ for each initial state \mathbf{s}_0^j
- 4. update policy $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ by imitating all $p_j(\mathbf{a}_t | \mathbf{s}_t)$





[parallelize sampling]
[parallelize dynamics]
[parallelize LQR]
[parallelize SGD]
3)

Yahya, Li, Kalakrishnan, Chebotar, L., '16

Model-based algorithms: parallel GPS



Real-world model-free deep RL: parallel NAF



NAF Architecture.

$$Q(\mathbf{x}, \mathbf{u} | \boldsymbol{\theta}^{Q}) = A(\mathbf{x}, \mathbf{u} | \boldsymbol{\theta}^{A}) + V(\mathbf{x} | \boldsymbol{\theta}^{V})$$
$$A(\mathbf{x}, \mathbf{u} | \boldsymbol{\theta}^{A}) = -\frac{1}{2} (\mathbf{u} - \boldsymbol{\mu}(\mathbf{x} | \boldsymbol{\theta}^{\mu}))^{T} \boldsymbol{P}(\mathbf{x} | \boldsymbol{\theta}^{P}) (\mathbf{u} - \boldsymbol{\mu}(\mathbf{x} | \boldsymbol{\theta}^{\mu}))$$





Gu*, Holly*, Lillicrap, L., '16

Simplest example: sample parallelism with off-policy algorithms

