

Meta-Learning

Outline

- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

Outline

- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

Background

- Machine VS Human
 - Machine have surpassed humans at many tasks.
 - However, Machine generally need far more data to reach the same level.
- Human can learn fast with a large amount of prior knowledge encoded in their brains and DNA!
- These prior knowledge are called “meta knowledge”!
- In fact, there is NO UNIFIED DIFINITION about WHAT IS META.....

Background

- **Meta-Learning——Learning to Learn.**
 - When you learn to play FPS, learn to play TCG, learn to play VR game, learn.....
----->You will know how to learn any game in a short time.
 - That is, You learned how to learn!
- The goal of meta-learning: Train a model on a variety of learning tasks, and solve new learning tasks using few training samples —— that is, **few-shot learning**.
- Requirement: Meta-learning should be general to the **task** and the form of **computation**. (similar to VIN, i.e. learn to plan)

Background

- This process can be viewed as building an **internal representation** that is broadly suitable for many tasks.
- Furthermore, it also means **maximizing the sensitivity** of the loss function of new tasks.

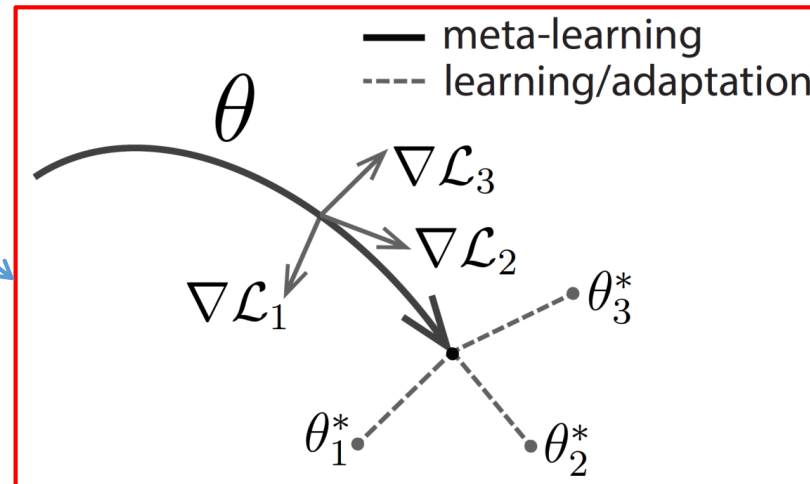
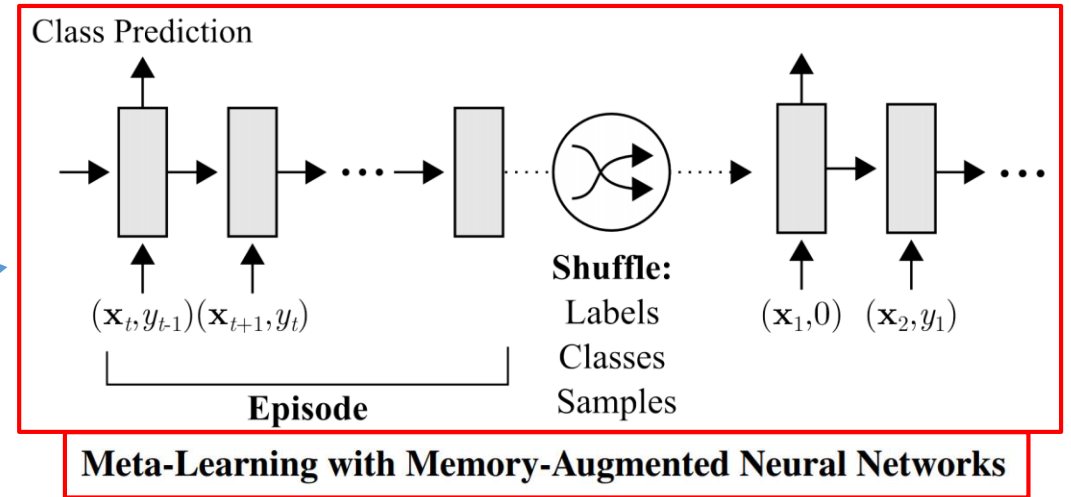
Background

- Three categories of methods:

- Recurrent Models
- Metric Learning

- Learning Optimizers

- Parameter Initialization**: learn the **initialization** of network and **fine-tune**(微调) on the new task, e.g. MAML.



Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks

Outline

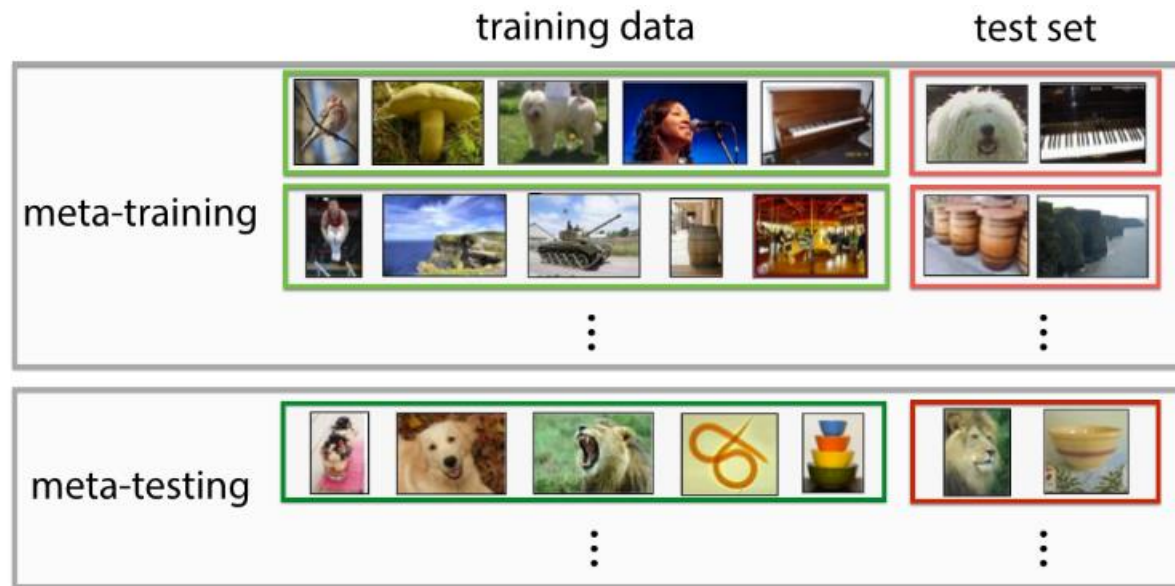
- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

Formulation

- Task: $\mathcal{T} = \{\mathcal{L}(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H, \mathbf{a}_H), q(\mathbf{x}_1), q(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{a}_t), H\}$
- Loss: $\mathcal{L}(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H, \mathbf{a}_H)$ (depend on specified task)
- Init Observation: $q(\mathbf{x}_1)$
- Transition Distribution: $q(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{a}_t)$
- Episode Length: H ($H=1$ for supervised learning)
- Task Distribution: $p(\mathcal{T})$
- **Goal:** Find a model to be able to adapt to the task distribution $p(\mathcal{T})$

Formulation

- How does meta-learning work to achieve this goal?
 - Meta-training and meta-testing.
 - Training set and test set.



Outline

- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

MAML

- Model-Agnostic Meta-Learning: a method to learn an initialization.
- Algorithm:

Meta
Training

- ① Initialization: meta-parameter θ
- ② Sample Tasks: $\mathcal{T}_i \sim p(\mathcal{T})$
- ③ Adaptation: compute adapted-parameter θ'_i using meta-parameter θ and training set, i.e. $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
- ④ Optimization: update θ using θ'_i and test set w.r.t. θ , i.e.

$$\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$$

why????

MAML

- Procedure ④ is an optimization problem:

$$\min_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i}) = \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})})$$

- That is to say, we want to find an optimal θ to make the performance of fine-tuning best! (i.e. **maximize the sensitivity**)
- However, the update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ involves **a gradient through a gradient**.
 - It requires additional computation(of course supported by some DL libraries)
 - Instead, we can only use a **first-order approximation**(FOMAML)

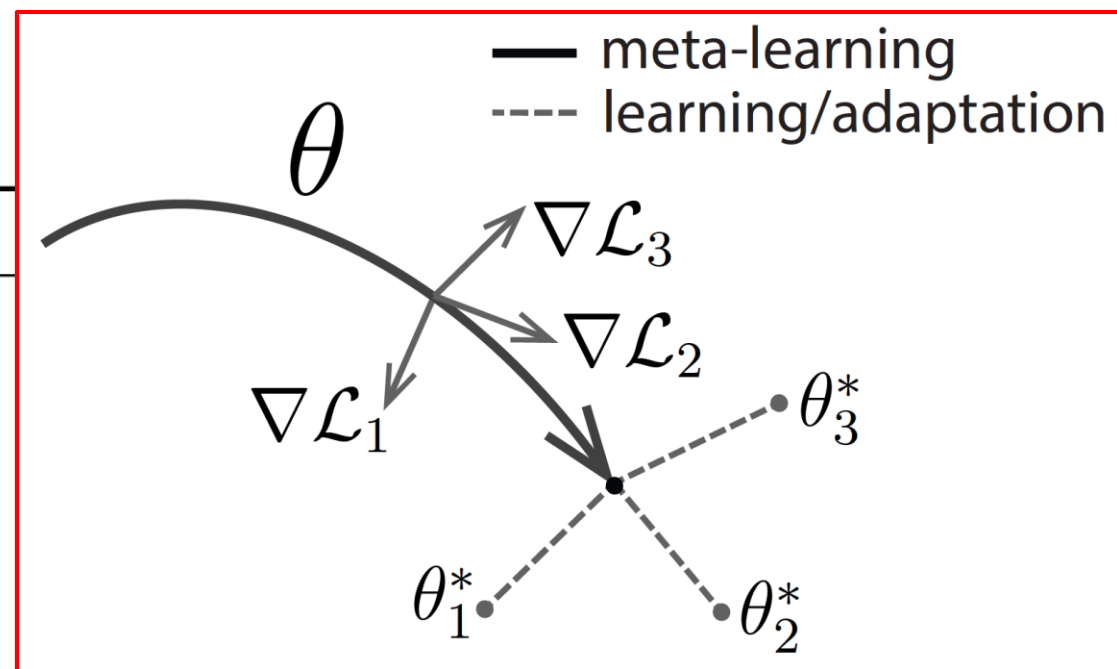
MAML

Algorithm 1 Model-Agnostic Meta-Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

- 1: randomly initialize θ
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ with respect to K examples
 - 6: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
 - 7: **end for**
 - 8: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
 - 9: **end while**
-



MAML

- With Algorithm 1, we can easily apply it to some different domains~~
- Supervised Learning
 - Regression
 - Classification
 -
- Reinforcement Learning
 - Policy Gradient
 - Actor Critic
 -

MAML

Algorithm 2 MAML for Few-Shot Supervised Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

- 1: randomly initialize θ
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: Sample K datapoints $\mathcal{D} = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$ from \mathcal{T}_i
 - 6: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation (2) or (3)
 - 7: Compute adapted parameters with gradient descent:
 $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
 - 8: Sample datapoints $\mathcal{D}'_i = \{\mathbf{x}^{(j)}, \mathbf{y}^{(j)}\}$ from \mathcal{T}_i for the meta-update
 - 9: **end for**
 - 10: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 2 or 3
 - 11: **end while**
-

MSE

$$\mathcal{L}_{\mathcal{T}_i}(f_{\phi}) = \sum_{\mathbf{x}^{(j)}, \mathbf{y}^{(j)} \sim \mathcal{T}_i} \|f_{\phi}(\mathbf{x}^{(j)}) - \mathbf{y}^{(j)}\|_2^2, \quad (2)$$

Cross-Entropy Loss

$$\begin{aligned} \mathcal{L}_{\mathcal{T}_i}(f_{\phi}) = \sum_{\mathbf{x}^{(j)}, \mathbf{y}^{(j)} \sim \mathcal{T}_i} & \mathbf{y}^{(j)} \log f_{\phi}(\mathbf{x}^{(j)}) \\ & + (1 - \mathbf{y}^{(j)}) \log(1 - f_{\phi}(\mathbf{x}^{(j)})) \end{aligned} \quad (3)$$

MAML

Algorithm 3 MAML for Reinforcement Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

- 1: randomly initialize θ
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: Sample K trajectories $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$ using f_θ
 - 6: Evaluate $\nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$ using \mathcal{D} and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 4
 - 7: Compute adapted parameters with gradient descent:
 $\theta'_i = \theta - \alpha \nabla_\theta \mathcal{L}_{\mathcal{T}_i}(f_\theta)$
 - 8: Sample trajectories $\mathcal{D}'_i = \{(\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H)\}$ using $f_{\theta'_i}$
 - 9: **end for**
 - 10: Update $\theta \leftarrow \theta - \beta \nabla_\theta \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$ using each \mathcal{D}'_i
 and $\mathcal{L}_{\mathcal{T}_i}$ in Equation 4
 - 11: **end while**
-

$$\pi_\theta(a | x_t)$$

$$\mathcal{L}_{\mathcal{T}_i}(f_\phi) = -\mathbb{E}_{\mathbf{x}_t, \mathbf{a}_t \sim f_\phi, q_{\mathcal{T}_i}} \left[\sum_{t=1}^H R_i(\mathbf{x}_t, \mathbf{a}_t) \right]. \quad (4)$$

MAML

- Benefits of MAML:
- 1.It is a **simple** framework.
- 2.It can be combined with **any model representation / architecture** and **any differentiable objective**. (MAML merely produces a parameter initialization)
- 3.Adaptation can be performed with **any amount of data** and **any number of gradient step**.
- As you can see, MAML is quite a **common framework**!

Outline

- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

MBMPO

- Model-Based Meta-Policy Optimization.
- Key Point:
 - Model-based approach with **ensemble**.
 - Meta-policy for fast adaptation.
- Benefits:
 - **Lower sample complexity**(model-based)
 - **Alleviate model bias to achieve better performance**(ensemble)
 - **Faster convergence**(meta-policy)

MBMPO

- Model Learning:

$$\min_{\phi_k} \frac{1}{|\mathcal{D}_k|} \sum_{(s_t, a_t, s_{t+1}) \in \mathcal{D}_k} \|s_{t+1} - \hat{f}_{\phi_k}(s_t, a_t)\|_2^2$$

models

- Meta-Reinforcement Learning:

$$\max_{\theta} \frac{1}{K} \sum_{k=0}^K J_k(\theta'_k) \quad \text{s.t.:} \quad \theta'_k = \theta + \alpha \nabla_{\theta} J_k(\theta)$$

meta-objective adaptation objective

- with the objective function: $J_k(\theta) = \mathbb{E}_{a_t \sim \pi_{\theta}(a_t | s_t)} \left[\sum_{t=0}^{H-1} r(s_t, a_t) \middle| s_{t+1} = \hat{f}_{\phi_k}(s_t, a_t) \right]$

MBMPO

Algorithm 1 MB-MPO

Require: Inner and outer step size α, β

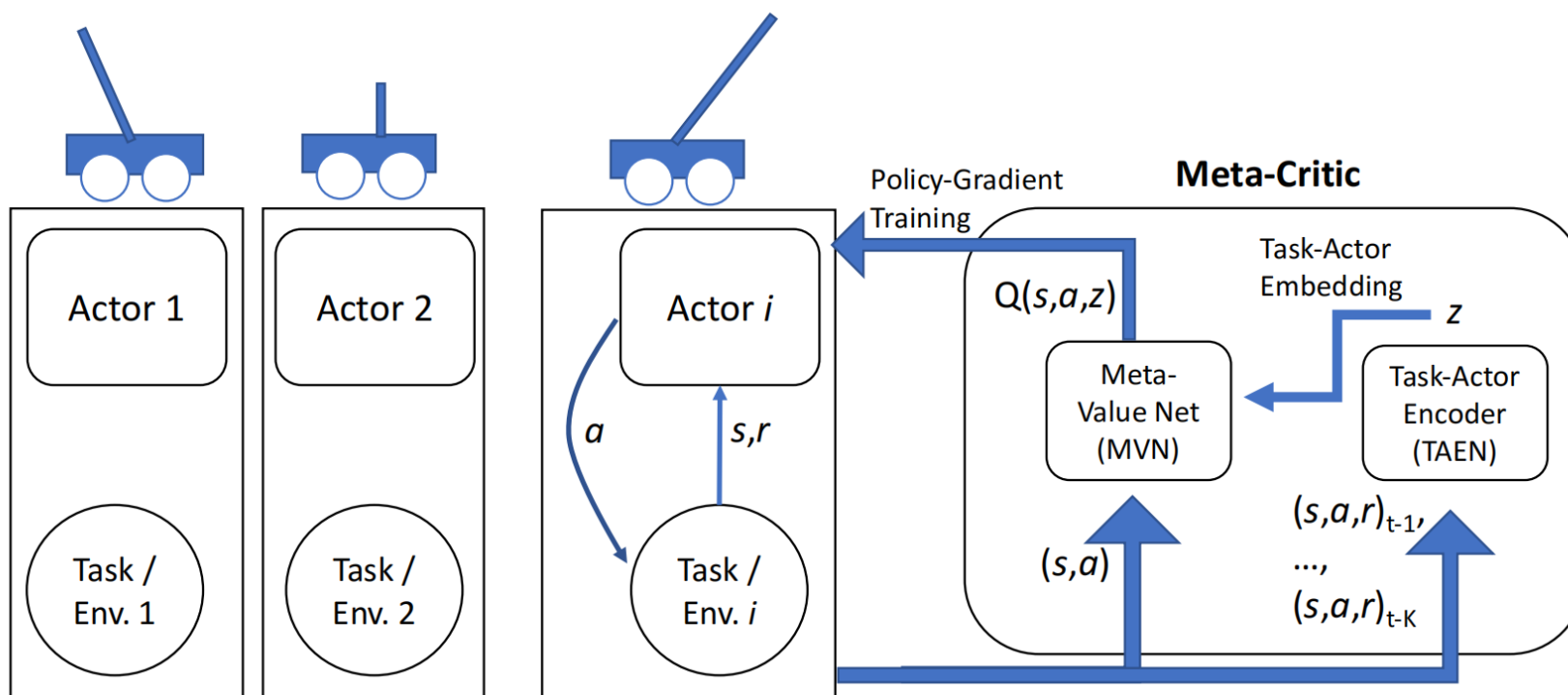
- 1: Initialize the policy π_{θ} , the models $\hat{f}_{\phi_1}, \hat{f}_{\phi_2}, \dots, \hat{f}_{\phi_K}$ and $\mathcal{D} \leftarrow \emptyset$
 - 2: **repeat**
 - 3: Sample trajectories from the real environment with the adapted policies $\pi_{\theta'_1}, \dots, \pi_{\theta'_K}$. Add them to \mathcal{D} .
 - 4: Train all models using \mathcal{D} .
 - 5: **for all** models \hat{f}_{ϕ_k} **do**
 - 6: Sample imaginary trajectories \mathcal{T}_k from \hat{f}_{ϕ_k} using π_{θ}
 - 7: Compute adapted parameters $\theta'_k = \theta + \alpha \nabla_{\theta} J_k(\theta)$ using trajectories \mathcal{T}_k
 - 8: Sample imaginary trajectories \mathcal{T}'_k from \hat{f}_{ϕ_k} using the adapted policy $\pi_{\theta'_k}$
 - 9: **end for**
 - 10: Update $\theta \rightarrow \theta - \beta \frac{1}{K} \sum_k \nabla_{\theta} J_k(\theta'_k)$ using the trajectories \mathcal{T}'_k
 - 11: **until** the policy performs well in the real environment
 - 12: **return** Optimal pre-update parameters θ^*
- PG method, e.g. TRPO, VPG,

Outline

- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

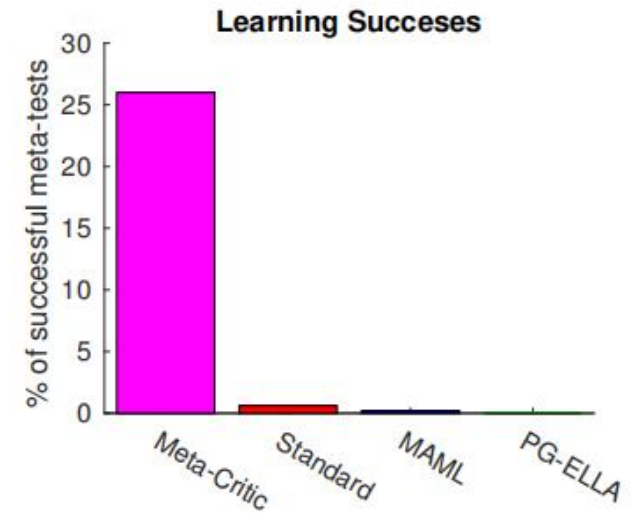
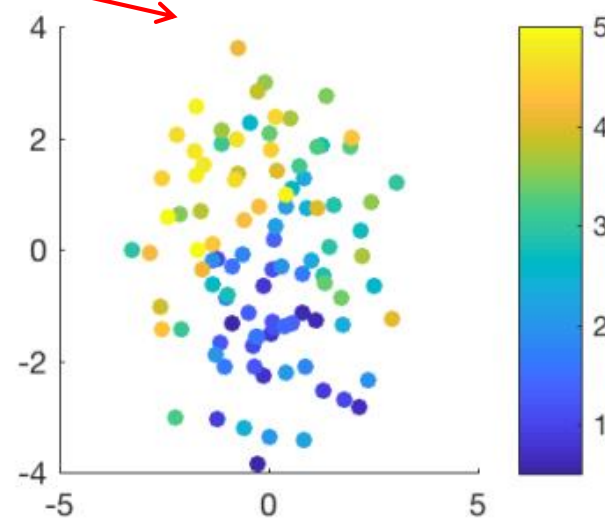
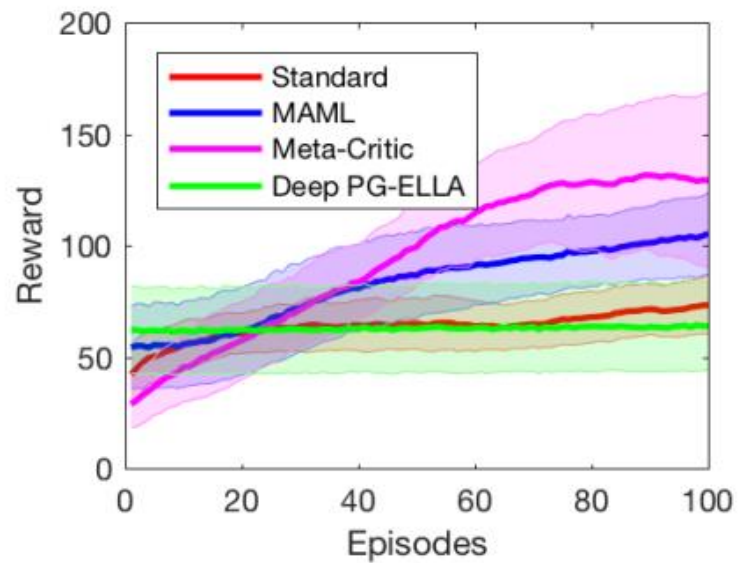
Meta-Critic-Networks

- Idea: Instead of learning a policy, Meta-Critic-Networks learn a **Meta Value Network(MVN)** and a **Task-Actor Encoder Network(TAEN)**.
- In their words, that is a “核心价值观”



Meta-Critic-Networks

- Experiment: Cartpole
- What does TAEN encode?



Meta-Critic-Networks

Algorithm 1: Meta-Learning Stage

Input: Task generator \mathcal{T}

Output: Trained task and value net

```
1 Init: task and value net;
2 for episode = 1 to max episode do
3   Generate  $M$  tasks from  $\mathcal{T}$ ;
4   Init  $M$  policy nets (actors);
5   for step = 1 to max steps do
6     Sample mini-batch of tasks;
7     foreach task in mini-batch do
8       Sample training data from task;
9       Train task-specific actor;
10    end
11    Train value network;
12    Train task network;
13  end
14 end
```

Meta-Critic-Networks

Algorithm 2: Meta-Testing Stage

Input: An unseen task

Input: Trained task and value nets

Output: Trained policy network

1 Init: one policy network (actor);

2 **for** *step = 1 to max step* **do**

3 | Sample train data from task;

4 | Train actor;

5 **end**

Outline

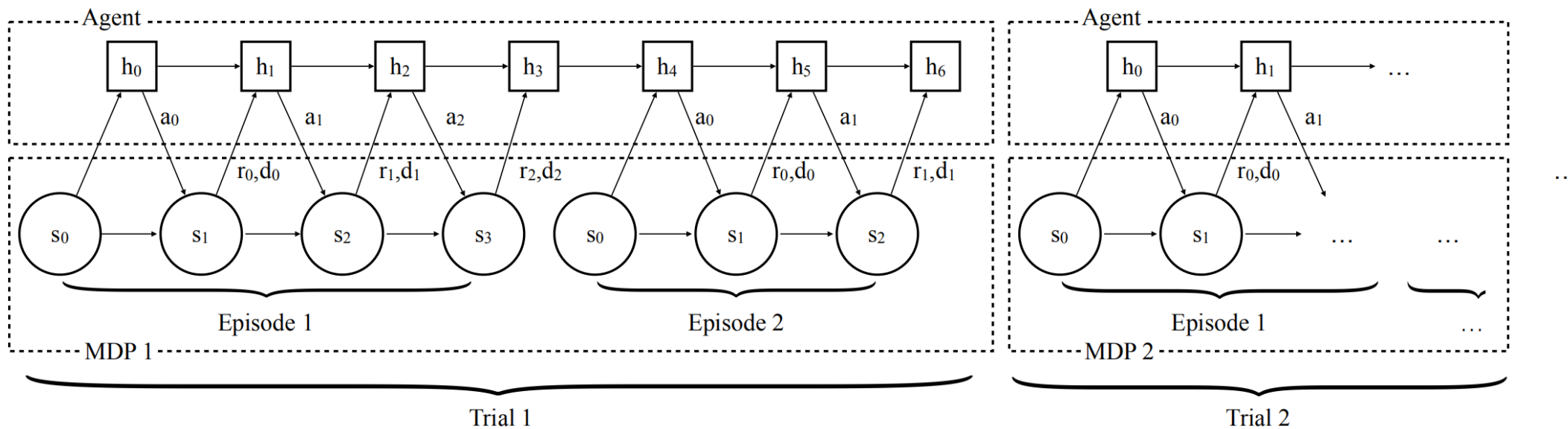
- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

RL²

- RL²: Casts learning an RL algorithm as a reinforcement learning problem
- Idea:
 - View the learning process of the agent itself as an objective.
 - Structure the agent as a **RNN**, which receives **past rewards, actions**, and **termination flags** as inputs in addition to the normally received **observations**.

RL²

- Illustration



- When MDP changes, the agent must act differently according to its belief over which MDP it is currently in. Hence, the agent is forced to integrate all the information and adapt its strategy continually.

Outline

- Background
- Formulation
- MAML
- MBMPO
- Meta-Critic-Networks
- RL^2
- Conclusion

Conclusion

- What is meta? What is meta-learning?
 - No unified definition.....
 - However, we still can make progress on it.
- MAML is a kind of common framework to produce a parameter initialization.
- Meta-Critic-Network produce a novel idea for meta-learning.
- RL² use RL to learn a RL process.
-

Reference

- <https://bair.berkeley.edu/blog/2017/07/18/learning-to-learn/>
- Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks
- Model-Based Reinforcement Learning via Meta-Policy Optimization
- On First-Order Meta-Learning Algorithms
- Learning to learn: Meta-critic networks for sample efficient learning
- RL²: FAST REINFORCEMENT LEARNING VIA SLOW REINFORCEMENT LEARNING