

THIRD-PERSON IMITATION LEARNING

Ren Yanjie

2021.08.05

2021 WDS summer seminar

Content

- Part 1: Introduction
- Part 2: GAN and GAIL
- Part 3: Third-Person Imitation Learning
- Part 4: Conclusion

Part 1: Introduction

- 1.1 Imitation Learning Concepts
- 1.2 Imitation Learning Methods

1.1 Imitation Learning Concepts

- Motivation:

- No Reward Function Available

- Setting:

- Input: Expert Demonstration $\{\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_K\}$ $\hat{\tau}_i = \langle s_1^i, a_1^i, s_2^i, a_2^i, \dots, s_{n_i+1}^i \rangle$
- State Function $p(\cdot|s_t, a_t)$ is available

- Goal

- To Get A Policy Function $\pi(a|s)$

- Main Methods:

- Behaviour Cloning (BC)
- Inverse Reinforcement Learning (IRL)
- **Generative Adversarial Imitation Learning (GAIL)**



1.2 Imitation Learning Methods

- input label

 - Behaviour cloning
 - $\hat{\tau} = \{(s_1, \hat{a}_1, s_2, \hat{a}_2, \dots)\} \xrightarrow{\text{green arrow}} \{(s_1, \hat{a}_1), (s_2, \hat{a}_2), \dots\}$

↑

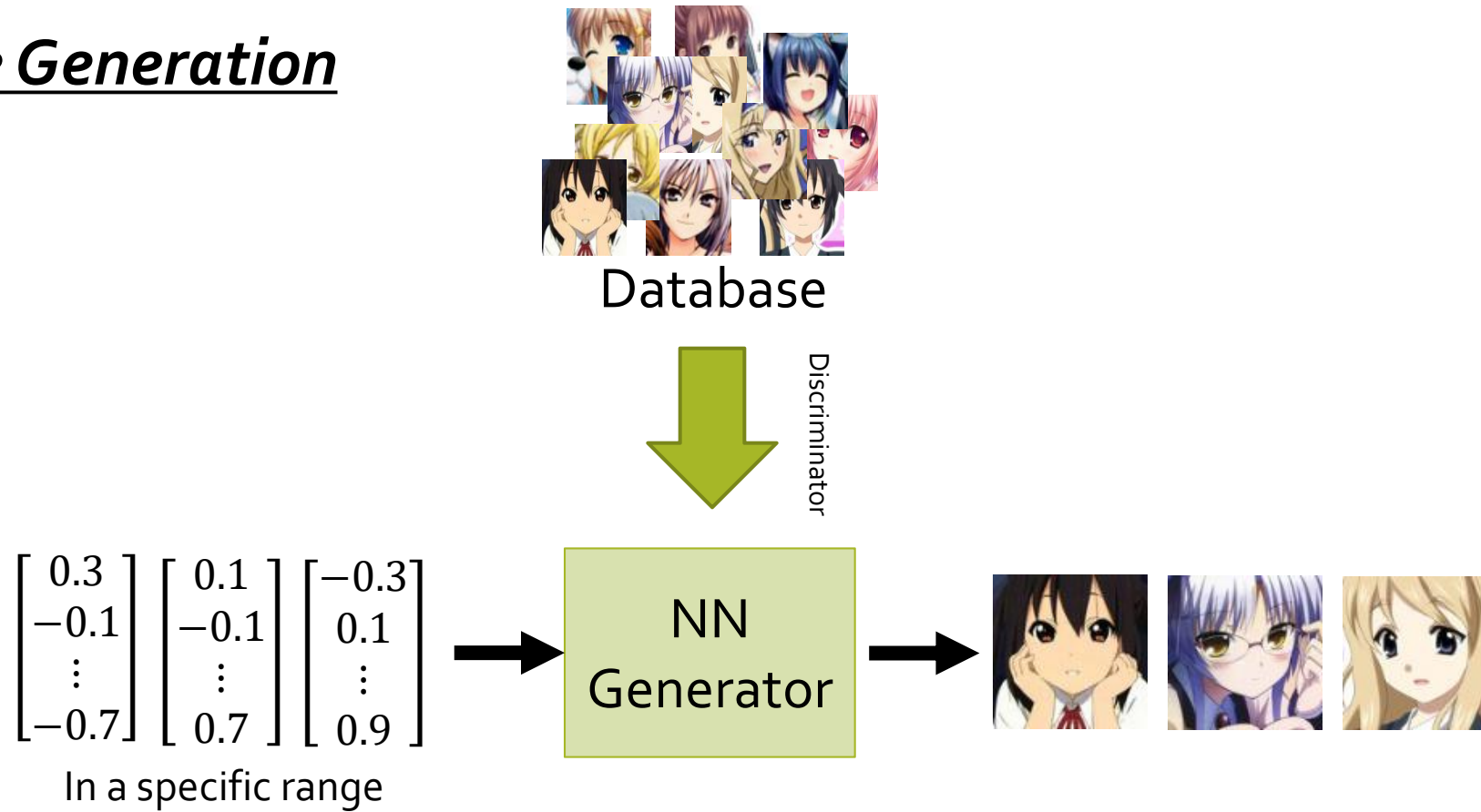
↑
 - Pros: intuitive, efficient, *off-line*
 - Cons: limit of train data, error-accumulate(step-data), **never better than expert**
 - Used to initialize model
 - Inverse Reinforcement Learning
 - Train The Reward Function ➔ Reinforcement Learning
 - Pros: solved error-accumulate problem
 - Cons: **time-consuming**, linear reward function restriction

Part 2: GAN and GAIL

- 2.1 GAN
- 2.2 GAIL
- 2.3 Pros And Cons

2.1 GAN

Image Generation



2.1 GAN

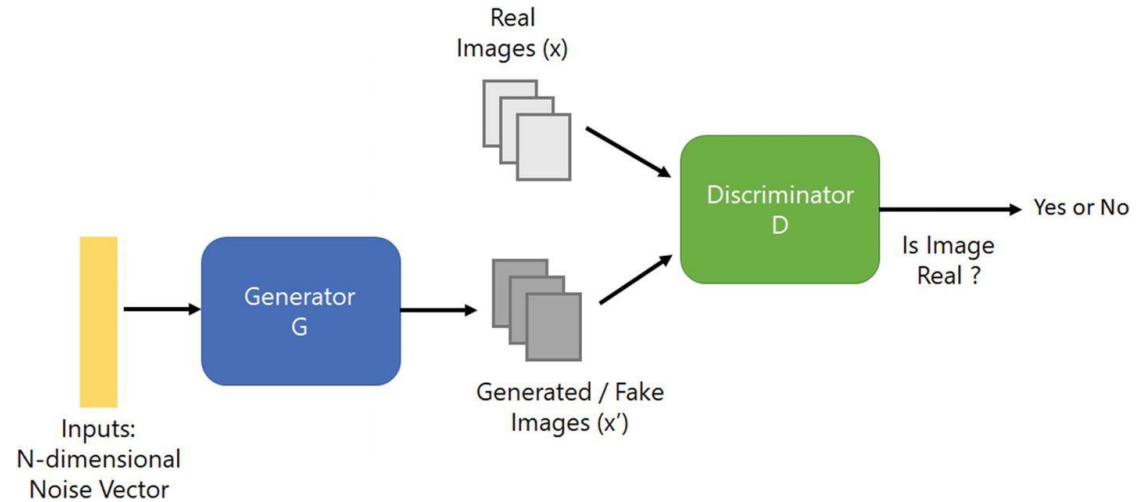
□ Basic Concepts:

□ Generator $G(s, \theta)$

- s : input vector in which $e \sim N(-1, 1)$ or $N(0,1)$
- θ : neural network argument
- $G(s_0, \theta)$: depends on the task

□ Discriminator $D(x, \varphi)$

- x : the output of generator or the sample data
- φ : the neural network argument
- $D(x, \varphi)$: a scalar which denotes the probability of sample data



2.1 GAN



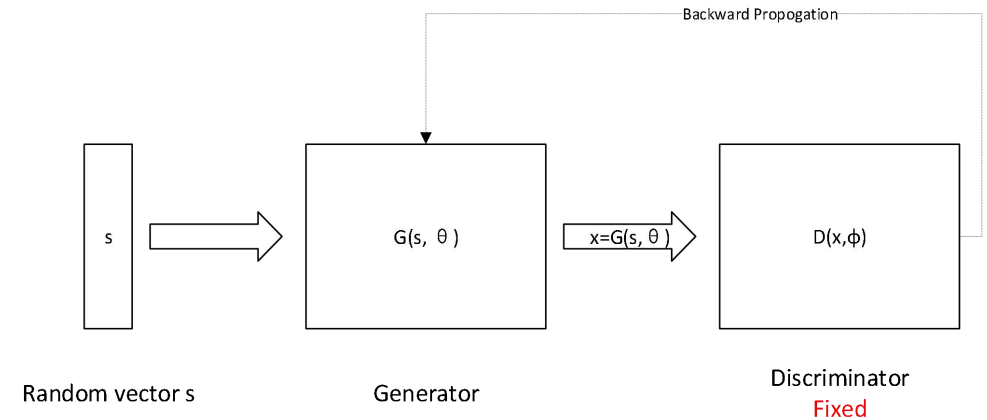
2.1 GAN

- ❑ Train the *generator*
 - ❑ Connect the generator $G(s, \theta)$ and discriminator $D(x, \phi)$
 - ❑ Fix the discriminator parameter ϕ
 - ❑ Define the loss function with cross entropy

$$E(s; \theta) = \ln \left[1 - \underbrace{D(x; \phi)}_{\text{越大越好}} \right]; \quad \text{s.t. } x = G(s; \theta).$$

- ❑ Update the θ with SGD or B-SGD

$$\theta \leftarrow \theta - \beta \cdot \nabla_{\theta} E(s; \theta).$$



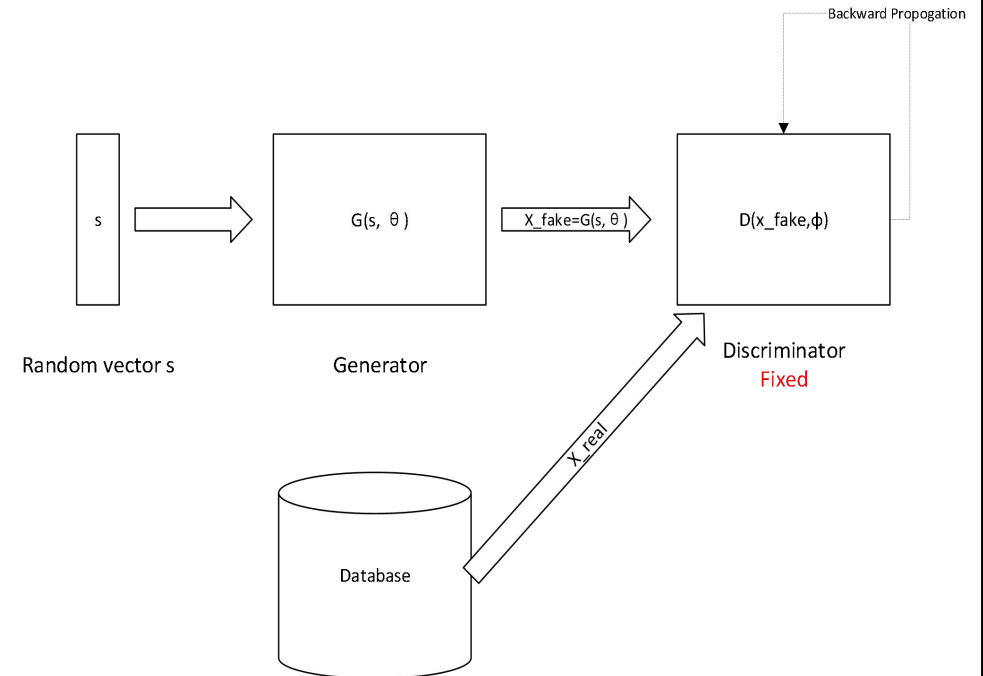
2.1 GAN

- ❑ Train the *discriminator*
 - ❑ Take a random sample from the database x^{real}
 - ❑ Get the $x^{\text{fake}} = G(s; \theta)$
 - ❑ Define the loss function with cross entropy

$$F(x^{\text{real}}, x^{\text{fake}}; \phi) = \ln \left[1 - \underbrace{D(x^{\text{real}}; \phi)}_{\text{越大越好}} \right] + \ln \underbrace{D(x^{\text{fake}}; \phi)}_{\text{越小越好}}.$$

- ❑ Update the ϕ with SGD or Mini-Batch SGD

$$\phi \leftarrow \phi - \eta \cdot \nabla_{\phi} F(x^{\text{real}}, x^{\text{fake}}; \phi).$$



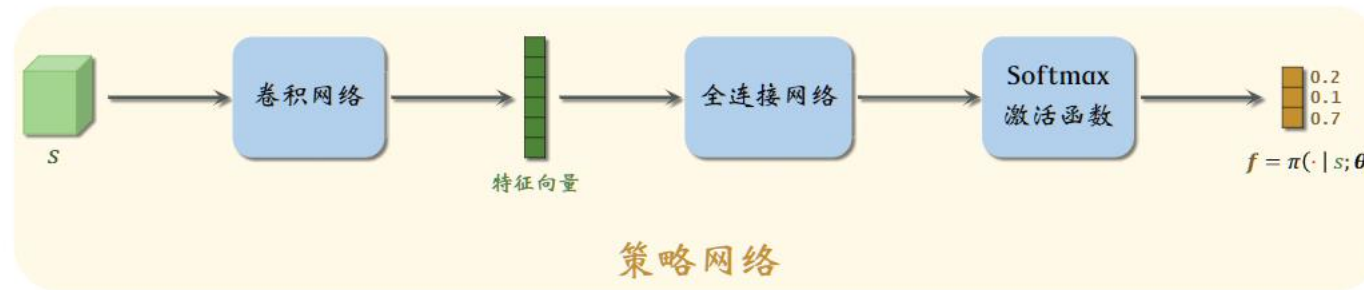
2.2 GAIL

❑ Expert Demonstration (Database)

- ❑ $\tau = [s_1, a_1, s_2, a_2, \dots, s_m, a_m]$.
- ❑ $\mathcal{X} = \{\tau^{(1)}, \tau^{(2)}, \dots, \tau^{(k)}\}$.

❑ Generator

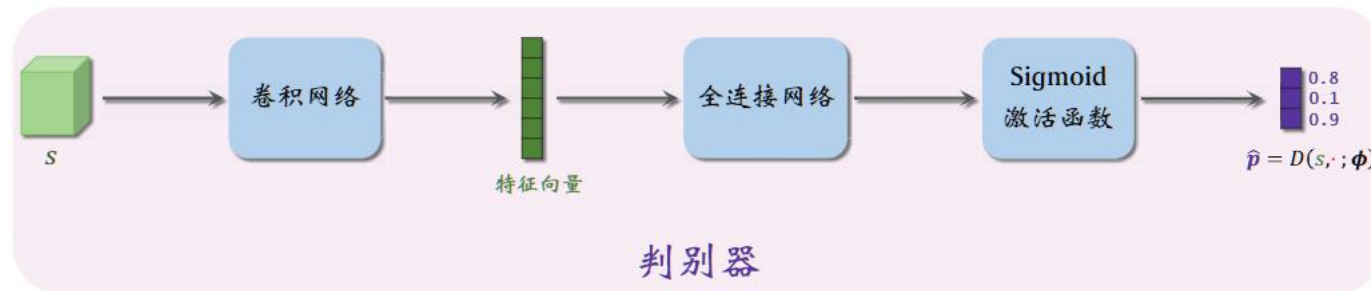
- ❑ $\pi(a|s; \theta) \dots\dots\dots G(s; \theta)$
- ❑ output $f = \pi(\cdot | s; \theta)$.
- ❑ given initial state s_1 , generator will get a trajectory $\tau = [s_1, a_1, s_2, a_2, \dots, s_n, a_n]$.



2.2 GAIL

□ Discriminator

- $\hat{p} = D(s, \cdot | \phi)$.
- $\hat{p}_a = D(s, a; \phi) \in (0, 1), \quad \forall a \in \mathcal{A}$.



2.2 GAIL: Algorithm

Algorithm 1 Generative adversarial imitation learning

- 1: **Input:** Expert trajectories $\tau_E \sim \pi_E$, initial policy and discriminator parameters θ_0, w_0
- 2: **for** $i = 0, 1, 2, \dots$ **do**
- 3: Sample trajectories $\tau_i \sim \pi_{\theta_i}$
- 4: Update the discriminator parameters from w_i to w_{i+1} with the gradient

$$\hat{\mathbb{E}}_{\tau_i} [\nabla_w \log(D_w(s, a))] + \hat{\mathbb{E}}_{\tau_E} [\nabla_w \log(1 - D_w(s, a))] \quad (17)$$

- 5: Take a policy step from θ_i to θ_{i+1} , using the TRPO rule with cost function $\log(D_{w_{i+1}}(s, a))$. Specifically, take a KL-constrained natural gradient step with

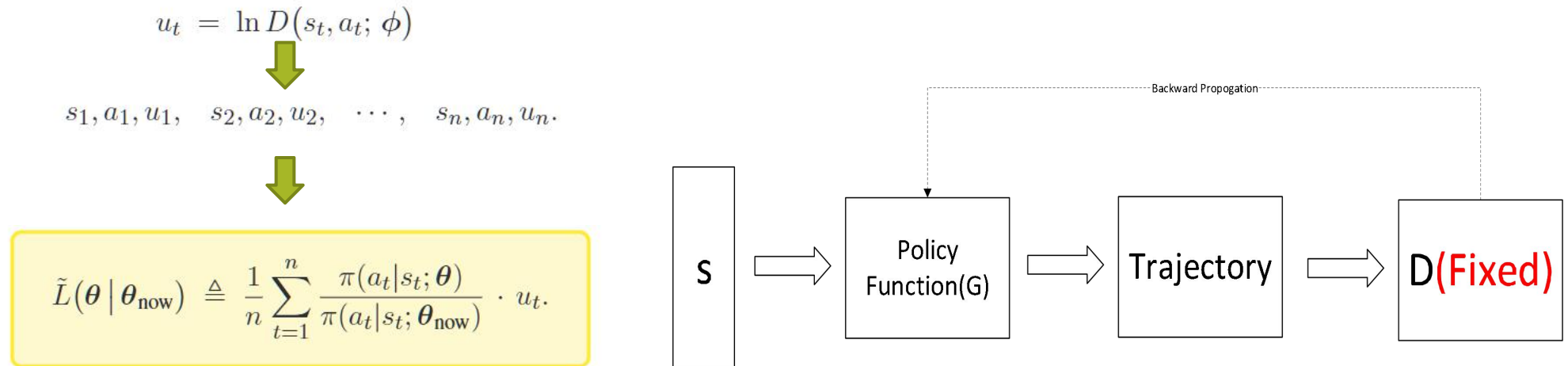
$$\hat{\mathbb{E}}_{\tau_i} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q(s, a)] - \lambda \nabla_{\theta} H(\pi_{\theta}), \quad (18)$$

where $Q(\bar{s}, \bar{a}) = \hat{\mathbb{E}}_{\tau_i} [\log(D_{w_{i+1}}(s, a)) \mid s_0 = \bar{s}, a_0 = \bar{a}]$

- 6: **end for**
-

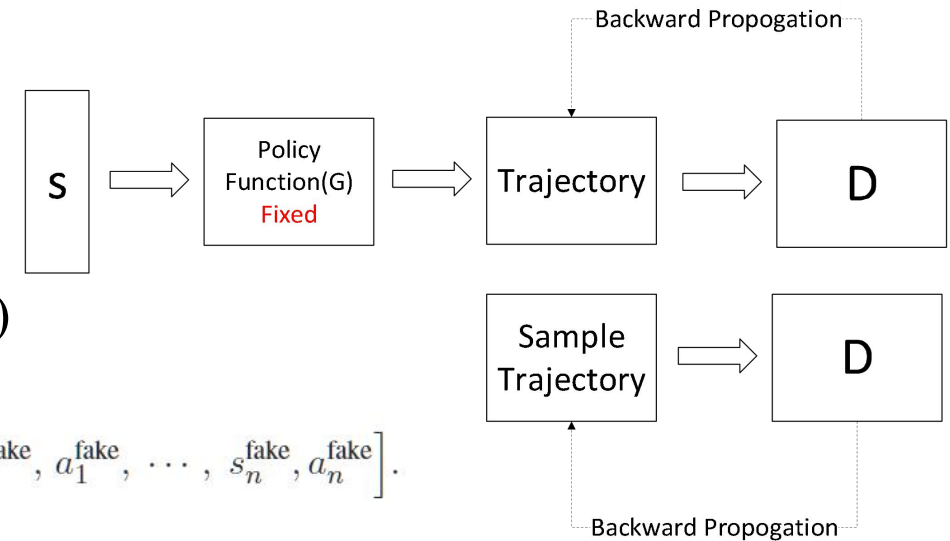
2.2 GAIL: Algorithm

- ❑ Train the generator (i.e. policy function $\pi(a|s; \theta_{\text{now}})$)
 - ❑ θ_{now} : current parameter of policy function
 - ❑ use current policy function to get a trajectory $\tau = [s_1, a_1, s_2, a_2, \dots, s_n, a_n]$.
 - ❑ $D(s_t, a_t; \phi)$ can output the quality of (s_t, a_t)



2.2 GAIL: Algorithm

- Train the discriminator (i.e. policy function $\pi(a|s; \theta_{\text{now}})$)
 - Random sample $\tau^{\text{real}} = [s_1^{\text{real}}, a_1^{\text{real}}, \dots, s_m^{\text{real}}, a_m^{\text{real}}]$.
 - Use current policy function to get a trajectory $\tau^{\text{fake}} = [s_1^{\text{fake}}, a_1^{\text{fake}}, \dots, s_n^{\text{fake}}, a_n^{\text{fake}}]$.
 - Define the loss function



$$F(\tau^{\text{real}}, \tau^{\text{fake}}; \phi) = \underbrace{\frac{1}{m} \sum_{t=1}^m \ln [1 - D(s_t^{\text{real}}, a_t^{\text{real}}; \phi)]}_{D \text{ 的输出越大, 这一项越小}} + \underbrace{\frac{1}{n} \sum_{t=1}^n \ln D(s_t^{\text{fake}}, a_t^{\text{fake}}; \phi)}_{D \text{ 的输出越小, 这一项越小}}.$$

2.3 GAIL: Pros and Cons

- ❑ Pros
 - ❑ Efficient Compared to IRL
 - ❑ Exact more features than IRL
- ❑ Cons
 - ❑ Modal collapse
 - ❑ Improved GAIL based on multiple modal
 - ❑ Improved GAIL based on
 - ❑ Low utilization of generated samples

Part 3: Third-Person Imitation Learning

- 3.1: Basic Concepts
- 3.2: DANN
- 3.3: TPIL

3.1 Third-Person

First Person

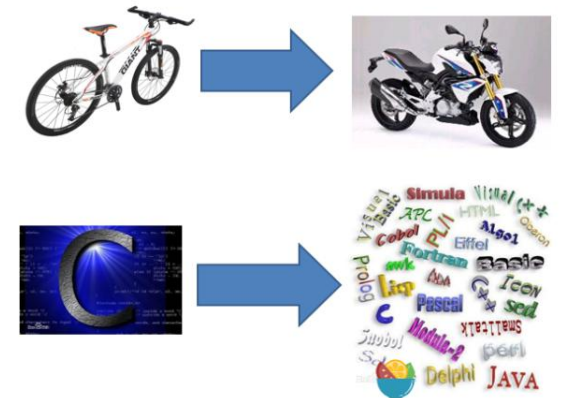


Third Person

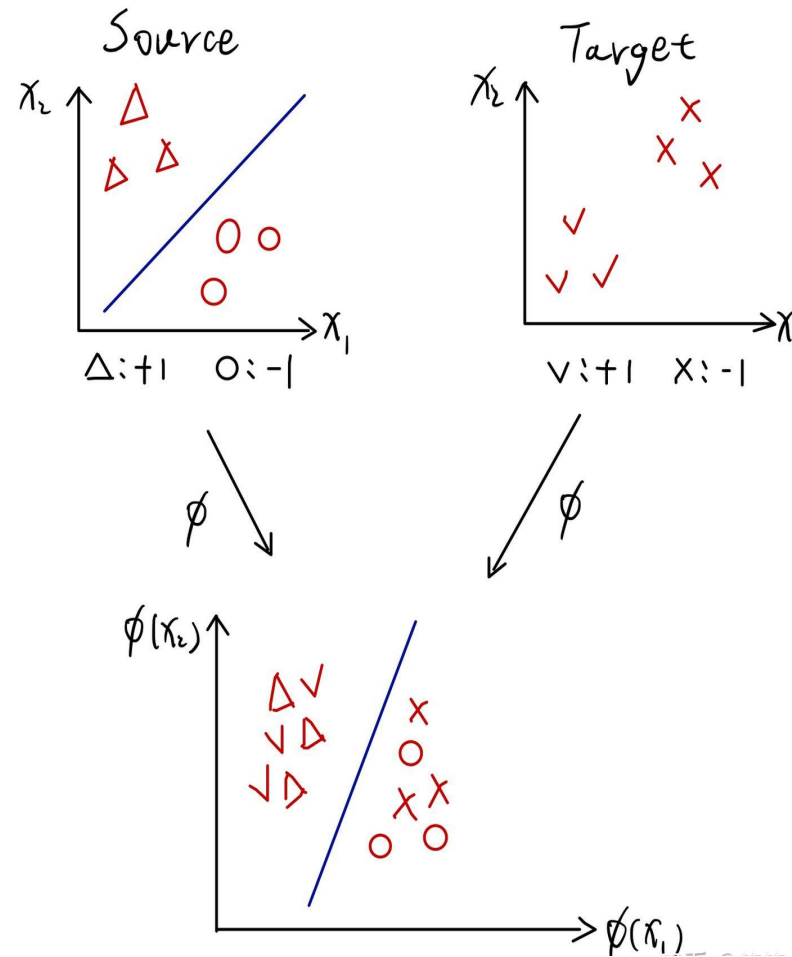


3.1 TPIL VS Transfer Learning

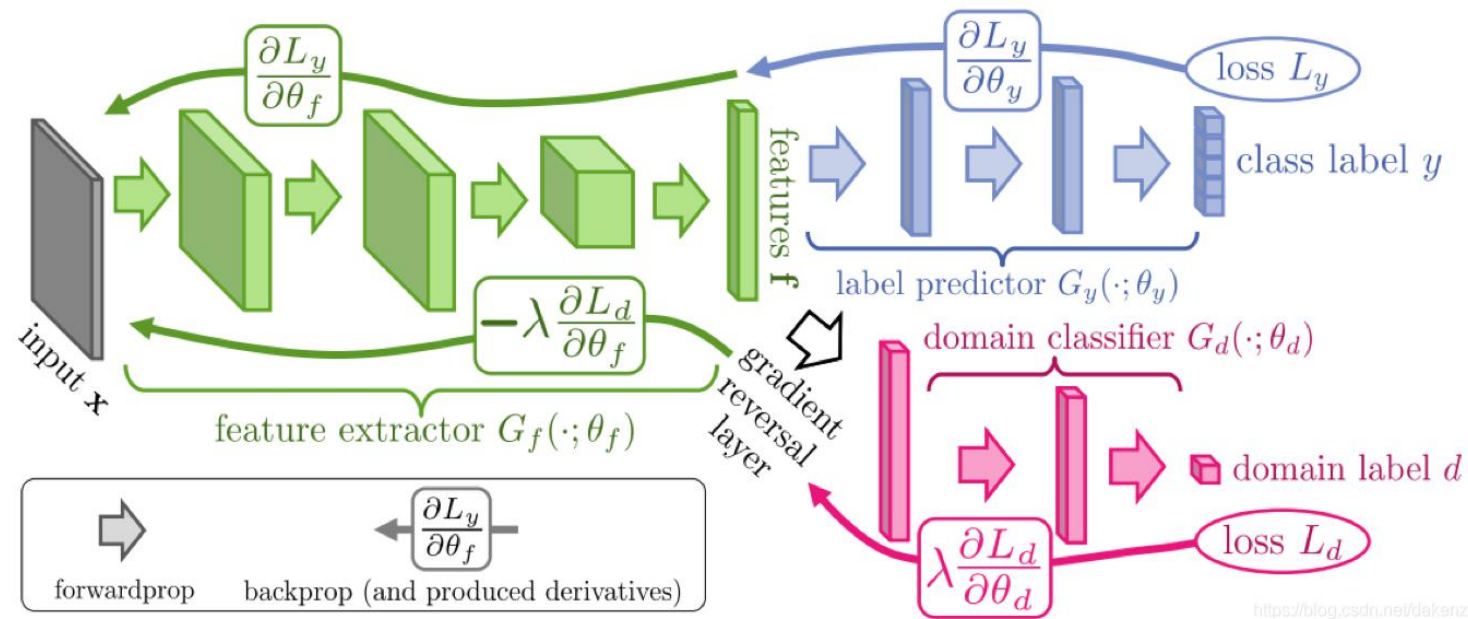
- Transfer Learning: Ability of a system to recognize and apply knowledge and skills learned in **previous domains/tasks to novel domains/tasks**
- Domain
- TPIL and Transfer Learning
 - Source Domain: Third person
 - Destion Domin: First Person



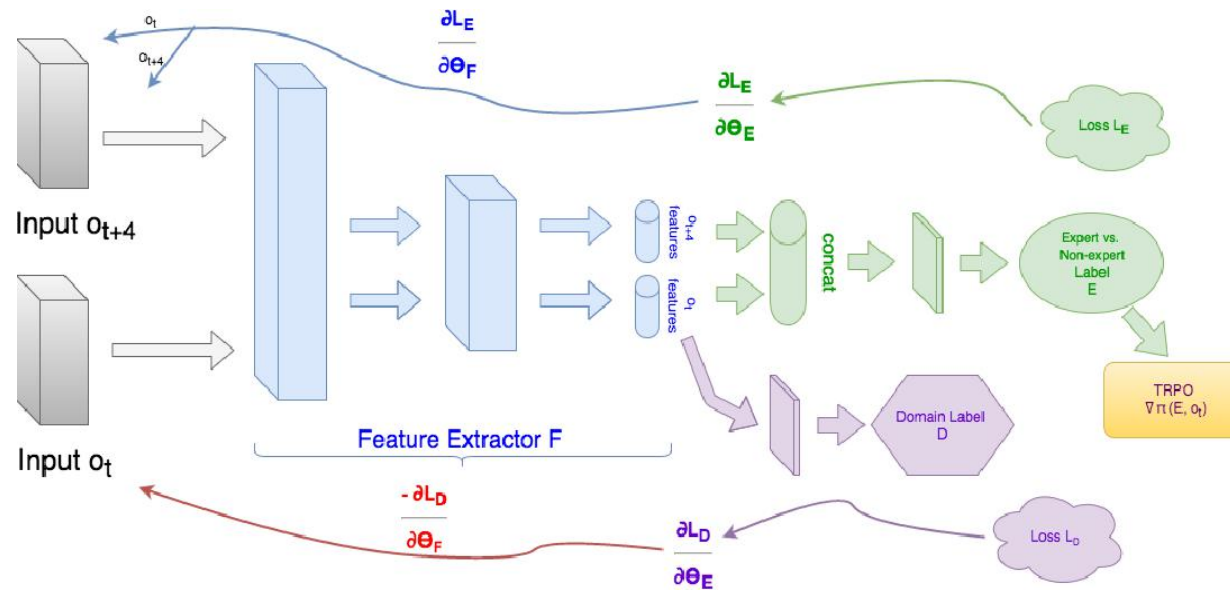
3.2 DANN



3.2 DANN



3.3 TPIL



3.3 TPIL

Algorithm 1 A third-person imitation learning algorithm.

```

1: Let CE be the standard cross entropy loss.
2: Let  $\mathcal{G}$  be a function that flips the gradient sign during backpropagation and acts as the identity map otherwise.
3: Initialize two domains,  $E$  and  $N$  for the expert and novice.
4: Initialize a memory bank  $\Omega$  of expert success and of failure in domain  $E$ . Each trajectory  $\omega \in \Omega$  comprises a rollout of images  $o = o_1, \dots, o_t, \dots, o_n$ , a class label  $c_\ell$ , and a domain label  $d_\ell$ .
5: Initialize  $\mathcal{D} = \mathcal{D}_F, \mathcal{D}_R, \mathcal{D}_D$ , a domain invariant discriminator.
6: Initialize a novice policy  $\pi_\theta$ .
7: Initialize numiters, the number of inner policy optimization iterations we wish to run.
8: for iter in numiters do
9:   Sample a set of successes and failures  $\omega_E$  from  $\Omega$ .
10:  Collect on policy samples  $\omega_N$ 
11:  Set  $\omega = \omega_E \cup \omega_N$ .
12:  Shuffle  $\omega$ 
13:  for  $o, c_\ell, d_\ell$  in  $\omega$  do
14:    for  $o_t$  in  $o$  do
15:       $\sigma_t = \mathcal{D}_F(o_t)$ 
16:       $\sigma_{t+4} = \mathcal{D}_F(o_{t+4})$ 
17:       $\mathcal{L}_R = CE(\mathcal{D}_R(\sigma_t, \sigma_{t+4}), c_\ell)$ 
18:       $\mathcal{L}_d = CE(\mathcal{D}_D(\mathcal{G}(\sigma_t)), d_\ell)$ 
19:       $\mathcal{L} = \lambda \cdot \mathcal{L}_d + \mathcal{L}_R$ 
20:      minimize  $\mathcal{L}$  with ADAM.
21:    end for
22:  end for
23:  Collect on policy samples  $\omega_N$  from  $\pi_\theta$ .
24:  for  $\omega$  in  $\omega_N$  do
25:    for  $\omega_t$  in  $\omega$  do
26:       $\sigma_t = \mathcal{D}_F(o_t)$ 
27:       $\sigma_{t+4} = \mathcal{D}_F(o_{t+4})$ 
28:       $\hat{c}_\ell = \mathcal{D}_R(\sigma_t, \sigma_{t+4})$ 
29:       $r = \hat{c}_\ell[0]$ , the probability that  $o_t, o_{t+4}$  were generated via expert rollouts.
30:      Use  $r$  to train  $\pi_\theta$  with via policy gradients (TRPO).
31:    end for
32:  end for
33: end for
34: return optimized policy  $\pi_\theta$ 

```



```

23: Collect on policy samples  $\omega_N$  from  $\pi_\theta$ .
24: for  $\omega$  in  $\omega_N$  do
25:   for  $\omega_t$  in  $\omega$  do
26:      $\sigma_t = \mathcal{D}_F(o_t)$ 
27:      $\sigma_{t+4} = \mathcal{D}_F(o_{t+4})$ 
28:      $\hat{c}_\ell = \mathcal{D}_R(\sigma_t, \sigma_{t+4})$ 
29:      $r = \hat{c}_\ell[0]$ , the probability that  $o_t, o_{t+4}$  were generated via expert rollouts.
30:     Use  $r$  to train  $\pi_\theta$  with via policy gradients (TRPO).
31:   end for
32: end for
33: end for
34: return optimized policy  $\pi_\theta$ 

```

3.3 TPIL

Algorithm 1 A third-person imitation learning algorithm.

```

1: Let CE be the standard cross entropy loss.
2: Let  $\mathcal{G}$  be a function that flips the gradient sign during backpropagation and acts as the identity map otherwise.
3: Initialize two domains,  $E$  and  $N$  for the expert and novice.
4: Initialize a memory bank  $\Omega$  of expert success and of failure in domain  $E$ . Each trajectory  $\omega \in \Omega$  comprises a rollout of images  $o = o_1, \dots, o_t, \dots, o_n$ , a class label  $c_\ell$ , and a domain label  $d_\ell$ .
5: Initialize  $\mathcal{D} = \mathcal{D}_F, \mathcal{D}_R, \mathcal{D}_D$ , a domain invariant discriminator.
6: Initialize a novice policy  $\pi_\theta$ .
7: Initialize numiters, the number of inner policy optimization iterations we wish to run.
8: for iter in numiters do
9:   Sample a set of successes and failures  $\omega_E$  from  $\Omega$ .
10:  Collect on policy samples  $\omega_N$ 
11:  Set  $\omega = \omega_E \cup \omega_N$ .
12:  Shuffle  $\omega$ 
13:  for  $o, c_\ell, d_\ell$  in  $\omega$  do
14:    for  $o_t$  in  $o$  do
15:       $\sigma_t = \mathcal{D}_F(o_t)$ 
16:       $\sigma_{t+4} = \mathcal{D}_F(o_{t+4})$ 
17:       $\mathcal{L}_R = CE(\mathcal{D}_R(\sigma_t, \sigma_{t+4}), c_\ell)$ 
18:       $\mathcal{L}_d = CE(\mathcal{D}_D(\mathcal{G}(\sigma_t)), d_\ell)$ 
19:       $\mathcal{L} = \lambda \cdot \mathcal{L}_d + \mathcal{L}_R$ 
20:      minimize  $\mathcal{L}$  with ADAM.
21:    end for
22:  end for
23:  Collect on policy samples  $\omega_N$  from  $\pi_\theta$ .
24:  for  $\omega$  in  $\omega_N$  do
25:    for  $\omega_t$  in  $\omega$  do
26:       $\sigma_t = \mathcal{D}_F(o_t)$ 
27:       $\sigma_{t+4} = \mathcal{D}_F(o_{t+4})$ 
28:       $\hat{c}_\ell = \mathcal{D}_R(\sigma_t, \sigma_{t+4})$ 
29:       $r = \hat{c}_\ell[0]$ , the probability that  $o_t, o_{t+4}$  were generated via expert rollouts.
30:      Use  $r$  to train  $\pi_\theta$  with via policy gradients (TRPO).
31:    end for
32:  end for
33: end for
34: return optimized policy  $\pi_\theta$ 

```



```

8: for iter in numiters do
9:   Sample a set of successes and failures  $\omega_E$  from  $\Omega$ .
10:  Collect on policy samples  $\omega_N$ 
11:  Set  $\omega = \omega_E \cup \omega_N$ .
12:  Shuffle  $\omega$ 
13:  for  $o, c_\ell, d_\ell$  in  $\omega$  do
14:    for  $o_t$  in  $o$  do
15:       $\sigma_t = \mathcal{D}_F(o_t)$ 
16:       $\sigma_{t+4} = \mathcal{D}_F(o_{t+4})$ 
17:       $\mathcal{L}_R = CE(\mathcal{D}_R(\sigma_t, \sigma_{t+4}), c_\ell)$ 
18:       $\mathcal{L}_d = CE(\mathcal{D}_D(\mathcal{G}(\sigma_t)), d_\ell)$ 
19:       $\mathcal{L} = \lambda \cdot \mathcal{L}_d + \mathcal{L}_R$ 
20:      minimize  $\mathcal{L}$  with ADAM.
21:    end for
22:  end for

```

Part 4: Conclusion

- 4.1: Imitation Learning
 - No Reward Function Available
 - Main methods: Behaviour Cloning/Inverse Reinforcement Learning
- 4.2: Generative Adversarial Imitation Learning
- 4.3: Third-person Learning
 - Extend of GAIL
 - DANN && TL && InfoGAN
- 4.4: Some open issues: Manually Setting Award Function / Techniques Innovation

Thanks!

Reference

- [1]林嘉豪,章宗长,姜冲,郝建业.基于生成对抗网络的模仿学习综述[J].计算机学报,2020,43(02):326-351.
- [2] Third-person imitation learning via image difference and variational discriminator bottleneck (student abstract version) . JIANG C,ZHANG Z Z,CHEN Z X,et al. Proceedings of the44th AAAI Conference on Artificial Intelligence . 2020
- [3] Third-person imitation learning. Stadie B C,Abbeel P,Sutskever I. Proceedings of the 5th International Conference on Learning Representations (ICLR) . 2017
- [4] Deep domain confusion:Maximizing for domain invariance. Tzeng E,Hoffman J,Zhang N,et al. . 2014
- [5] Generative adversarial imitation learning. Ho J,Ermon S. Neural Information Processing Systems . 2016
- [6] Domain-Adversarial Training of Neural Networks
- [7] Learning and Transferring Mid-Level Image Representations using Convolutional Neural Networks