

1) Imitation Learning from observations by minimizing inverse dynamics disagreement

### General Notes

- minimizing the discrepancy in the inverse dynamics models of expert and agent.
- model free
- Upper bound: negative entropy of state action occupancy measure
- Different occupancy measures for MDP :-
  - i) State Action
  - ii) State Transition
  - iii) Joint Occupancy
- Agent and expert share same dynamics.
- Without access of dynamics model and expert guidance.

2) To Follow or Not to Follow

### General Notes

- learns to follow demonstrations by aligning the timescale and stripping infeasible parts.
- Hierarchy of policies :-
  - i) Meta → Sub Goal
  - ii) Low-Level → How to get to
- Currents observations and demonstrations are collected by different agents in different environments, thus embeddings are used.



### 3) IL by State Only Distribution Matching

#### General Notes

- AIL methods are often unstable & lack a reliable convergence estimator.
- This paper proposes a non adversarial LFO approach.
- Training objective minimizes the KLD
- Estimates the expert state transition distribution using normalizing flows (trained offline)
- To match transition distributions of policy and expert exactly, the state-next-state distribution of policy is expanded into policy entropy, forward dynamics and inverse action model

### 4) Decoupled Policy Optimization

#### General Notes

- Decouples Policy as a high level state planner and an inverse dynamics model.
- Uses embedded decoupled policy gradient with GAN.
- Takes a supervised approach, but has GAT

### 5) Differentiable Physics

- Incorporates the differential physics simulator as a physics prior into its computational graph for policy learning.
- ILD unrolls dynamics by sampling actions from a parameterized policy and minimizing distance.
- Selects learning objectives for each state dynamically.

## 6) Raw Video via Context Translation

### General Notes

- Context translation + deep RL
- Deals with domain shift

## 7) Self Supervised Adversarial IL

- By incorporating a discriminator
  - i) Disposes manual intervention
  - ii) Guiding function approximation based on the state transition of expert trajectories
  - iii) Avoids a no action performance

## 8) MOBILE

- Uses the distribution matching IL framework.



# ① Imitation from Observation

$$p(o_t | s_t, w); p(s_{t+1} | s_t, a_t, w); p(a_t | s_t, w)$$

$w \rightarrow$  context

Further challenge: Domain shift

~~2~~ 2 problems:  $\rightarrow$  i) what info from observations to track in own context.

ii) Actions to track demonstrated observations

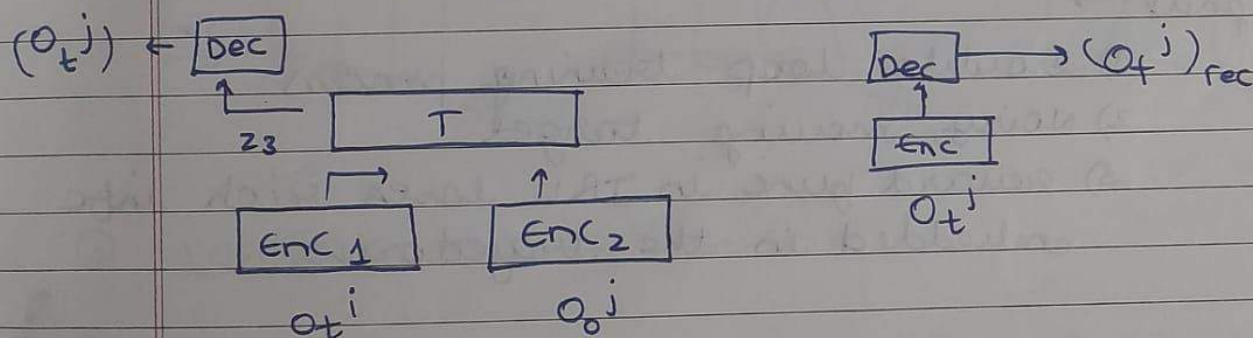
$$D_i = [o_0^i, o_1^i, \dots, o_T^i] \rightarrow \text{source context}$$

$$D_j = [o_0^j, o_1^j, \dots, o_T^j] \rightarrow \text{target context}$$

$\rightarrow$  "The model first learns to output observations in  $D_j$  conditioned on  $D_i$  and first obs  $o_0^j$  and context  $w_j$ ."

$\rightarrow$  Assume  $D_i, D_j$  are aligned in time.

$$M(o_t^i, o_0^j) = (o_T^j)_{\text{trans}}$$



$$L_{\text{Trans}} = \| (\hat{o}_t^j)_{\text{trans}} - o_t^j \|_2^2$$

$$L_{\text{rec}} = \| \text{Dec}(\text{Enc}(o_t^j)) - o_t^j \|_2^2$$

$$L_{\text{align}} = \| z_3 - \text{Enc}(o_t^j) \|_2^2$$

$$L = L_t + \lambda_1 L_r + \lambda_2 L_a$$

hyperparams



## → Reward Func's for Tracking Features

$$\hat{R}_{\text{feat}}(o_t^l) = -\| \text{enc}_1(o_t^l) - \frac{1}{n} \sum^i \text{enc}_1(o_t^i, o_0^l) \|_2^2$$

- ④ Distribution of observations fed into  $\text{enc}_1$  may not be same during policy learning and training. We need a reward that directly penalizes the policy for exp observations that differ from translated obs using  $M$ :-

$$\hat{R}_{\text{imp}}(o_t^l) = -\| o_t^l - \frac{1}{n} \sum^i M(o_t^i, o_0^l) \|_2^2$$

## ② Imitation Learning as State Matching via Differentiable Physics

- AIL and IRL have three main limitations due to learning an additional intermediate signal: →

- 1) Double loop training process
- 2) Noisy, moving target
- 3) reward pure in IRL loses rich info embedded in the trajectories.

$$x_{t+1} = x_t + \Delta t v_t, \quad v_{t+1} = v_t + \frac{f \Delta t}{m}$$

→  $\Pi_\theta$  is used to unroll → trajectory  $T_0 \xrightarrow{S_{0:H}} A_{0:H}$   
 $T(S_{t+1} | S_t, a_t)$  is our env. dynamics.

$$S_{1:H} = G(S_0, a_{0:H})$$

$$\arg \min_{\theta} \sum_{S \sim T_\theta} \sum_{t=0}^T (g_t - s_t)^2$$



$\alpha_{\text{triplet}} \Rightarrow$  enforces encodings sim of act frames

$\alpha_{\text{ae}} \Rightarrow$  permits enc to be dec to image

classmate

Date

Page

$$s = g_{\theta}(0), s_p = g_{\theta}(0_p)$$

$$\mathcal{L}_{\text{ae}}(\theta, \gamma) = \|0 - g_{\gamma}(g_{\theta}(0))\|^2$$

Finally,

$$\{R, a, B\} \rightarrow \{L, ab\}$$

$$s = g(\theta) = [g_{\theta_1}(v_1), g_{\theta_2}(v_2)]$$

$$\mathcal{L}(\theta, \gamma) = \mathcal{L}_{\text{frame}}(\theta) + \mathcal{L}_{\text{triplet}}(\theta) + \mathcal{L}_{\text{ae}}(\theta, \gamma)$$

### ③ Sequence encoding

i)  $L_2$  (DPC)

$$D = \{0^i\}_{i \leq N} = \{0^i_{0:T}\}_{1 \leq i \leq N} \rightarrow \text{dataset of } D \text{ indices}$$

$$\begin{aligned} \hat{s}_{t+1} &= d_{\phi}(f_w(s_1 \dots s_t)) \\ \hat{s}_{t+k} &= d_{\phi}(f_w(s_1 \dots \hat{s}_{t+k-1})) \end{aligned}$$

$$L_2(\theta, w, \phi) = -\mathbb{E}_{0:T} \left[ \frac{1}{K} \sum_{1 \leq k \leq K} L_{t+k} \right]$$

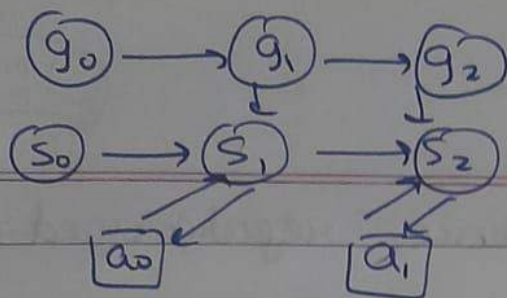
$$L_t = \log \frac{h(\hat{s}_t, s_t)}{h(\hat{s}_t, s_t) + \sum_{t \neq t'} h(\hat{s}_t, s_{t'})}$$

where  $h(a, b) = \exp\left(\frac{a^T b}{\|a\| \|b\|}\right)$   $T \Rightarrow \text{temp}$

- $L_0 \Rightarrow$  bootstrapped expert behaviour loss, which encourages separation between expert and non expert trajectories

$$L_0(\theta, w) = -\mathbb{E} \log \frac{h(z, z_p)}{h(z, z_p) + \sum_i h(z, z_{n,i})}$$

$$z = f_w(g_{\theta}(0_0) \dots g_{\theta}(0_T))$$



classmate

Date

Page

We use Chamfer- $\alpha$  Loss  $\rightarrow$  Deviation Loss + Coverage Loss

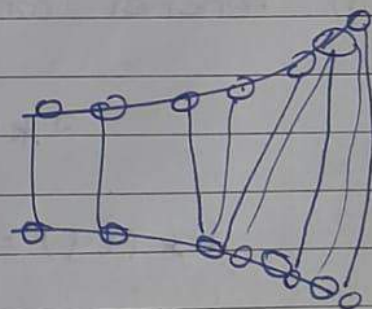
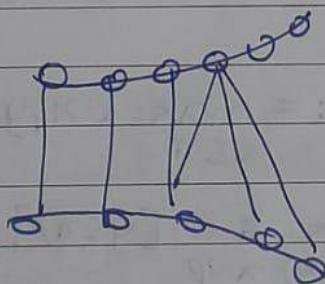
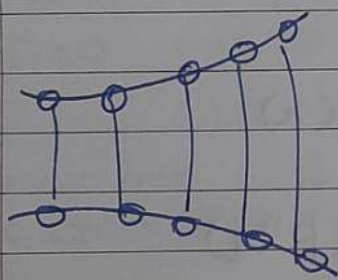
$$L_d = \frac{1}{|T_0|} \sum_{s_t \in T_0} \min_{g \in T_{exp}} \|g - s_t\|_2^2$$

$\Rightarrow$  But, this may lead to state collapse.

eg. Consider a large number of states close to small subset of  $T_{exp}$ .

$$L_g = \frac{1}{|T_{exp}|} \sum_{g_t \in T_{exp}} \min_{s \in T_{\pi}} \|g_t - s\|_2^2$$

$$L_{chf-\alpha} = L_d + \alpha L_g$$



③ LOBSDice  $\rightarrow$  (verify for SOLL)

④ Versatile offline Imitation from observations...

High Level :-  $\min_{\pi} D_{KL}(d^{\pi}(s) || d^E(s))$

- Naively optimizing the state occupancy distribution would result in an actor critic style IL algorithm.



$$= \min_{\pi} - E_{(s,a) \sim d^{\pi}} [E^*(s,a) \log \pi(a|s)]$$

where  $E^*(s,a) = \frac{d^*(s,a)}{d^{\pi}(s,a)}$

- ⑤ Imitation from observation with Bootstrapped CL  
 → 2 main approaches for IFO are - reward method  
 - adversarial methods

Advantages / Motivation behind this approach.

- (i) To learn and estimate at each timestep the system state using agent's visual obs.
- (ii) Identify behaviour induced by agent and make sure it serves same purpose as expert

## ① Agent Training

- 2 phases → Alignment / Interactive

### (i) Alignment

$$O \sim p(D_e), O \sim p(D_a)$$

$f_w \Rightarrow$  sequence encoding function.

$g_o \Rightarrow$  image encoding function

### (ii) Interactive

PRQ-v2 ← → RL task where reward = distance( $O_e, O_a$ )

### ④ Sample an expert episode

$$O_e \sim p(D_e)$$

$$d(O_{0:t}, O_{e 0:t}) = \|f_w(g_o(O_{0:t})) - f_w(g_o(O_{e 0:t}))\|$$

$$r_t = -d(O_{0:t}, O_{e 0:t})$$

continue to train  $f_w, g_o$  on new trajectories

## ② Image encoding

$$\alpha_{\text{triplet}}(\theta) = \|s - s_p\|^2 + \max(p - \|s - s_n\|^2, 0)$$



- To bring sequences from the same distribution close together we use  $L_0$ .  
 For  $\theta$  belonging to one sequence,  
 $\theta_p \Rightarrow$  same distribution  
 $\theta_{n,i} \Rightarrow$  different distribution

$$\alpha_{seq}(\theta, w, \phi) = L_2(\theta, w, \phi) + L_0(\theta, w)$$

## ⑥ Policy Contrastive Imitation Learning

- Problem :- AIL  $\rightarrow$  Low quality of discriminator representation, since it is trained via binary classification.
- use contrastive learning loss
- Push the experts representation together and pull the agents away from them.

### ⇒ CPIL

- Common heuristics dictate that  $x_p$  is just an augmentation to  $x$ , but this is not strong consideration ~~point~~ for AIL, since we need to discern good/bad behaviour.
- Good and bad states may look similar or the difference can be minor.
- This causes a diff b/w positive sample and another sample; overwhelms the diff b/w good/bad.

$$\alpha = \mathbb{E} \left[ \phi(x_0)^T \phi(x_p) + \log \left( \exp \phi(x_0)^T \phi(x_p) + \sum_i \exp \phi(x_0)^T \phi(x_i) \right) \right]$$

$$r(x) = \phi(x)^T \mathbb{E}_{x \in \mathcal{D}} \phi(x)$$



## ⑦ versatile skill control via self-supervised

- Given a large dataset of unlabelled motion clips with diverse behaviours, extracting and learning individual sensible skills can be challenging.
- One promising attempt for unsupervised skill discovery is based on maximizing the discriminability of skills represented by latent variables on which policy is conditioned.
- Policy training signals are often constructed from variational approximations of the mutual information between latent variables and traversed state histories using a learned skill discriminator.

EXISTING  $\Rightarrow$  1) collect disparate reference motions and learn individual skills separately.  
2) stack training discriminators over separate tasks

- APPROACH:  $\Rightarrow$  i) imitation discriminator ( $d_\psi$ )  
ii) skill discriminator ( $q_\phi$ )

$$O_t^I = (O_{t-H^I+1}^I \dots O_t^I)$$

$$S_t^I = (S_{t-H^I+1}^I \dots S_t^I)$$

- 2) Imitation discriminator Formulation
- $$\mathbb{E}_{d^M} [(d_\psi(O^I) - 1)^2] + \mathbb{E}_{d^M} [(d_\psi(f^I(s)) + 1)^2]$$

$$r^I = \max[0, 1 - 0.25(d_\psi(f^I(s)) - 1)^2]$$

$$r \in [0, 1]$$



- Introducing  $f$ -divergence regularized state matching objective.
- Using dual optimality solution to formulate a weighted regression policy objective that amounts to BC of optimal policy.

Prelim  $d^\pi(s, a) := (1 - \gamma) \sum \gamma^t \Pr(s_t = s, a_t = a)$

$$d^\pi(s, a) = (1 - \gamma) \mu_0(s) \pi(a|s) + \gamma \cdot T_*^\pi d^\pi(s, a)$$

$$T_*^\pi d^\pi(s, a) := \pi(a|s) \sum T(s|\tilde{s}, \tilde{a}) d(\tilde{s}, \tilde{a})$$

i)  $f$ -Divergence

$$D_f(p||q) = \mathbb{E}_{x \sim q} \left[ f\left(\frac{p(x)}{q(x)}\right) \right]$$

ii) Fenchel conj

$$f_*(y) := \max_{x \in \mathbb{R}} \langle x, y \rangle - f(x)$$

$$D_* f(y) = \max_{p \in \Delta(X)} \mathbb{E}_{x \sim p} [y(x)] - D_f(p||q)$$

$$= \mathbb{E}_{x \sim q} [f_*(y(x))]$$

Algo 1 Assumption-1  $d^q(s) > 0$  when  $d^E(s) > 0$   
 $\hookrightarrow$  coverage

Assumption-2

$$D_{KL}(d^\pi(s)||d^E(s)) \leq$$

$$\mathbb{E}_{s \sim d^\pi} \left[ \log \left( \frac{d^0(s)}{d^E(s)} \right) \right] + D_{KL}(d^\pi(s, a)||d^0(s, a))$$

Any f-divergence  $D_f \geq D_{KL}$ ,

f-divergence  
regularized state  
occupancy match-  
ing objective

$$D_{KL}(d^\pi(s) || d^E(s)) \leq \mathbb{E}_{s \sim d^\pi} \left[ \log \left( \frac{d^0(s)}{d^E(s)} \right) \right] + D_f(d^\pi(s, a) || d^0(s, a))$$

## ① Discriminator Training

$$\min_c \mathbb{E}_{s \sim d^E} [\log c(s)] + \mathbb{E}_{s \sim d^0} [\log (1 - c(s))]$$

$$c^*(s) = \frac{d^0(s)}{d^0(s) + d^E(s)}$$

## ② Dual Value Function Training

$$\max_{d(s, a) \geq 0} \mathbb{E}_{s \sim d(s, a)} [R(s)] - D_f(d || d^0)$$

$$\begin{aligned} \max_{d(s, a) \geq 0} \min_{v(s) \geq 0} \mathbb{E}_{s \sim d} [R(s)] - D_f(d || d^0) + \sum_s v(s) (1 - \gamma) u_0(s) \\ + \gamma T_* d(s) - \sum_a d(s, a) \\ \sum_s v(s) \cdot T_* d(s) = \sum_a d(s, a) \cdot TV(s, a) \end{aligned}$$

$$\Rightarrow \min_{v(s) \geq 0} \max_{d(s, a) \geq 0} (1 - \gamma) \mathbb{E}_{s \sim u_0} [v(s)] + \mathbb{E}_{(s, a) \sim d} [(R(s) + \gamma TV(s, a) - v(s))] - D_f(d(s, a) || d^0(s, a))$$

$$\Rightarrow \min_{v(s) \geq 0} (1 - \gamma) \mathbb{E}_{s \sim u_0} [v(s)] + \mathbb{E}_{(s, a) \sim d^0} [f_*(R(s) + \gamma TV(s, a) - v(s))]$$

## ③ Weighted Regression Policy Training

$$\min_{\pi} - \mathbb{E}_{(s, a) \sim d^*} [\log \pi(a|s)]$$





$$J(\pi_\theta) = \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \mu^{\pi_\theta}} [r(s_t, a_t) + \gamma V(\pi_\theta(\cdot|s_t))]$$

• Look at soft actor critic method.

# Method

$$\mu^{\pi_\theta}(s_T, \dots, s_0) = P(s_0) \prod_{i=0 \dots T-1} \mu^{\pi_\theta}(s_{i+1}|s_i)$$

$$\mu^{\epsilon}(s_T, \dots, s_0) = P(s_0) \prod_{i=0 \dots T-1} \mu^{\epsilon}(s_{i+1}|s_i)$$

Goal is to match state-only trajectory distribution  $\mu^{\pi_\theta}$  induced by the policy with the state only expert trajectory distribution  $\mu^{\epsilon}$  by minimizing the KL between them.

$$J_{\text{soft-TDM}} = D_{KL}(\mu^{\pi_\theta} \parallel \mu^{\epsilon})$$

$$= \mathbb{E}_{(s_T, \dots, s_0) \sim \mu^{\pi_\theta}} [\log \mu^{\pi_\theta} - \log \mu^{\epsilon}]$$

$$= \sum_{i=0 \dots T-1} \mathbb{E}_{(s_{i+1}, s_i) \sim \mu^{\pi_\theta}} [\log \mu^{\pi_\theta}(s_{i+1}|s_i) - \log \mu^{\epsilon}(s_{i+1}|s_i)]$$

$$\mu^{\pi_\theta}(s_{i+1}|s_i) = \frac{P(s_{i+1}|a_i, s_i) \pi_\theta(a_i|s_i)}{\pi_\theta(a_i|s_i, s_{i+1})}$$

$$\min \sum_{i=0 \dots T-1} \mathbb{E}_{(s_i, a_i, s_{i+1}) \sim \pi_\theta} [\log P(s_{i+1}|a_i, s_i) + \log \pi_\theta(a_i|s_i)]$$

$$= \log \pi_\theta(a_i|s_{i+1}, s_i) - \log \mu^{\epsilon}(s_{i+1}|s_i)$$

$$r(a_i, s_i) := \mathbb{E}_{s_{i+1} \sim P(s_{i+1}|s_i, a_i)} [-\log P(s_{i+1}|a_i, s_i)]$$

$$+ \log \pi_\theta(a_i|s_i, s_{i+1}) + \log \mu^{\epsilon}(s_{i+1}|s_i)]$$

$$\min D_{KL}(\mu^{\pi_\theta} \parallel \mu^{\epsilon})$$

$$= \max \sum_{i=0 \dots T-1} \mathbb{E}_{(a_i, s_i) \sim \pi_\theta} [-\log \pi_\theta(a_i|s_i) + r(a_i, s_i)]$$

$$= \max \sum_{i=0 \dots T-1} \mathbb{E}_{(a_i, s_i) \sim \pi_\theta} [r(a_i, s_i) + H(\pi_\theta(\cdot|s_i))]$$

2) Reverse dynamics model how to code it?

• Bayes formulation

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$P(s_{i+1}|s_i) = \frac{P(s_i|s_{i+1})P(s_{i+1})}{P(s_i)}$$

$$P(s_{i+1}|s_i) = P(s_i, a_i|s_i)$$

$$\pi(a_i|s_i)P(s_{i+1}|s_i, a_i)$$

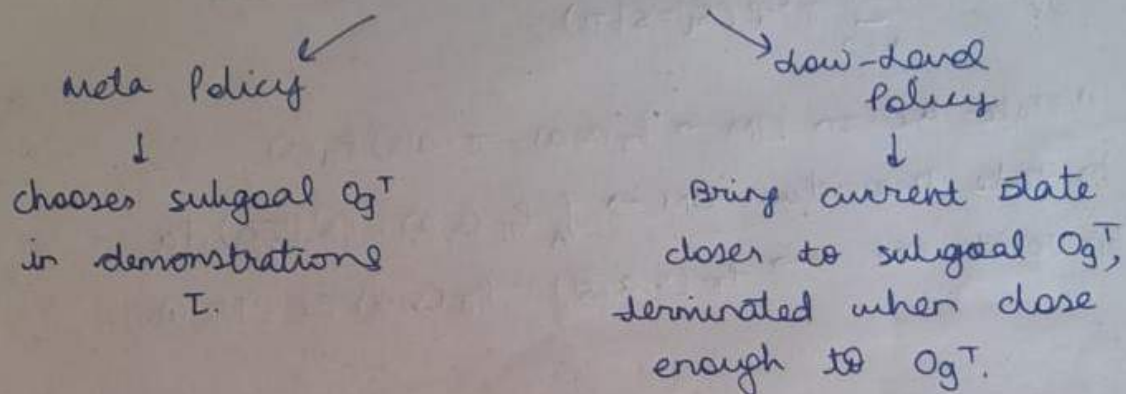
$$P(s_i)$$



## Silo Paper.

Naively following a demonstration does not work when demonstration consists a state that is unreachable by the learner agent.

- Uses a hierarchical RL framework



$$\pi_{\text{meta}}(g | o_t, \tau; \theta)$$

$o_t \Rightarrow$  current observation

$\tau \Rightarrow \{o_1^T, o_2^T, o_3^T, \dots, o_T^T\} \Rightarrow$  demonstration to

$g \in [1, T] \Rightarrow$  demonstration state <sup>initiate</sup>  
index of sub goal.

Once meta policy chooses  $o_g^T$  as a subgoal, the goal conditioned low-level policy generates an action  $a_t \sim \pi_{\text{low}}(a | o_t, o_g^T; \phi) \Rightarrow$  generates a rollout

until  $|o_g^T - o_{t+\Delta}| < \epsilon$ .

Reward for meta policy  $\Rightarrow 1/0$

⊕ uses embeddings to condense into common vector space.

## SAIL (Self Supervised Adversarial Imitation Learning)

### 4 Models →

- i)  $M$   $P(a|s_t, s_{t+1})$  → Inverse Dynamics Model
- ii)  $\pi_\theta$  → Policy Model
- iii)  $G_\pi$  → Generative Model
- iv)  $D$  → Discrimination Model

### Overview of Algorithm

- i) Random weights initialization
- ii) Use  $\pi_\theta$  to collect samples as  $(s_t, a, s_{t+1})$   $(I_t)$
- iii) Train  $M$ , and predict pseudo labels  $\hat{a}$  for  $T_\pi$
- iv) Train  $T_\theta$  using BE Approach
- v) Also updates  $G_\pi$  during training
- vi) Use  $T_\pi$  (updated) to create new samples for  $M$  ~~to train~~ <sup>action  $T_\pi$</sup>
- vii) Append to  $I_t$  all samples  $D$  cannot differentiate between  $T_\pi$  and  $G_\pi(I_{\text{new}})$

### Goal Aware Function

$$\min_{M, \pi, D} \text{SAIL}(M, \pi, G, D) = \mathbb{E} [\log(D(I_t))] \\ \min_{\pi, G} \mathbb{E} [\log(\pi(G, G_\pi))] \\ + \mathbb{E} [\log(1 - D(G(G_\pi, \pi(G_\pi))))]$$

$$\min_{\pi, G} \mathbb{E} [\log(1 - D(G(G_\pi, \pi(G_\pi))))]$$

NOTE: Beginning bad samples appended to  $I_t$ .

To avoid overfitting early, we deploy buffer

$$R(s_t, \dots, s_{t+n}) \sim T^e \cup T^\pi$$

### Generative Model

↳ "used at slow learning time".  
 ↳ selecting samples appended to  $I_t$

→ SAIL allows  $G_\pi$  to update  $\pi_\theta$  to create actions to correct condition  $T_\pi$  to update  $\pi_\theta$ , ~~create actions~~ generate correct state transitions, equal to above observed.  $B_t$

→  $G_\pi$  becomes a forward dynamics model.

Problems

- ↳  $N$  may not always be correct
- ↳  $N$  may not be stuck at some local minima.

$$L_G = -\frac{1}{N} \left[ \sum_{i=1}^N s_{i+1} \cdot \log(G(s_i, \pi(s_i))) \right]$$

ex)  $D$  dividing updating weights →  $\pi_\theta$  can direct-learned speed at where it deviates.



② the task is formulated as  $\rightarrow$

$$D = \{x_1^d, x_2^d, \dots\}$$

$\hookrightarrow$  i) At state  $x_0$ , take action  $\pi(x_0, x_1^d, \theta) \rightarrow x_0'$

If goal Recognizer( $x_0', x_1^d$ ) is High,  
 $x_0 = x_0'$ , take action.

~~~~~ while loop. till low.

② Learning GSP is entropy loss  $\rightarrow L(a_t, \tilde{a}_t)$   
 $\swarrow$   
 $p(a_t) \rightarrow$  actual  $\rightarrow$  predicted action dis

$\hookrightarrow$  Problem  $\rightarrow$  multimodality

③ Instead of penalizing actions, penalize closeness on next observed state from both distributions.

$$\rightarrow \min \|x_{t+1} - \tilde{x}_{t+1}\|^2 + \lambda \|x_{t+1} - \hat{x}_{t+1}\|^2 + L(a_t, \hat{a}_t)$$

$\hookrightarrow$  we need a good forward model 'f'.

# State Alignment Based Imitation Learning

\* 4 pointers

- i) Inverse State Based Model ✓
- ii) Deviation Correction (B-VAE)
- iii) Global Alignment
- iv) Regularized Policy Update

⇒ Local alignment → Policy Prior

Global Alignment → Rewards

⇒ VAE to produce next state

$\phi \Rightarrow$  discriminator

$$W = \mathbb{E}_{S \sim T_e} [\phi(s)] - \mathbb{E}_{S \sim T} [\phi(s)]$$

$$r(s_i, s_{i+1}) = \frac{1}{T} (\phi(s_{i+1}) - \mathbb{E}_{S \sim T_e} \phi(s))$$



## State Only Imitation Learning for Dexterous Manipulation

- ① similar to BCO, but joint iterative training.  
②

$$a_t' = h_\phi(s_t, s_{t+1}), L_{inv} = \|a_t' - a_t\|^2 \quad \nearrow \text{to train inverse model.}$$

- ② optimize  $J_{\text{soil}} = g + \lambda_0 d_1^k + \sum \nabla_\theta \log \pi_\theta(a_t | s)$

## adversarial imitation learning from state only demonstrations

① experts Data  $\rightarrow$  critic; imitator's Data  $\rightarrow$  actor

$\hookrightarrow D_\theta, \pi_\phi$ .

Take action according to  $\pi_\phi, \tau$ .  $\rightarrow$  actual trajectory

$$D_\theta \doteq -(\mathbb{E}_\tau [\log(D_\theta(s, s'))]) + \mathbb{E}_{\tau_e} [\log(1 - D_\theta(s, s'))].$$

$\pi_\phi \rightarrow$  TRPO updates.



## BCO

- ① <sup>Learn</sup> ~~the~~ an Inverse Dynamics Model  
by first letting agent ~~at~~ explore pre-demonstration. ( $\pi_\phi$ )
- ② we have 'D' which contains numerous trajectories.  
We use our inverse dynamics model on the agent specific part to infer action  $\tilde{a}_i$ .
- ③ ~~now we solve~~ cast this problem as BC.  $\Rightarrow$  ~~get~~ with  $(s_i, \tilde{a}_i)$  pairs.  
find  $\theta$  st probability of  $(\tilde{a}_i, \text{act} | s_i)$  is max.