SIMPLIFYING MODEL-BASED RL: LEARNING REPRESENTATIONS, LATENT-SPACE MODELS, AND POLICIES WITH ONE OBJECTIVE

Raj Ghugare^{1,2}Homanga Bharadhwaj²Benjamin Eysenbach²Sergey Levine³Ruslan Salakhutdinov²¹VNIT Nagpur²Carnegie Mellon University³UC Berkeleyraj19@students.vnit.ac.in, hbharadh@cs.cmu.edu, beysenba@cs.cmu.edu



0 Preliminaries - Reinforcement learning

Find policy $\pi(a_t|s_t)$ that maximises:

$$\max_{\pi} \mathbb{E}_{\underbrace{s_{t+1} \sim p\left(\cdot \mid s_{t}, a_{t}\right)}_{\text{environment}}, \underbrace{a_{t} \sim \pi\left(\cdot \mid s_{t}\right)}_{\text{policy}} \left[\underbrace{(1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} r\left(s_{t}, a_{t}\right)}_{R(\tau)}\right]$$

- Actions: *a*_t
- Reward function: $r(s_t, a_t) \ge 0$

(1)

0 Preliminaries - Model-based reinforcement learning

Find policy $\pi(a_t|s_t)$ that maximises:

$$\max_{\pi} \mathbb{E}_{\underbrace{s_{t+1} \sim p_{\theta}\left(\cdot \mid s_{t}, a_{t}\right)}_{\text{dynamics}}, \underbrace{a_{t} \sim \pi\left(\cdot \mid s_{t}\right)}_{\text{policy}} \left[\underbrace{(1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} r\left(s_{t}, a_{t}\right)}_{R(\tau)}\right]$$

- Actions: *a*_t
- Reward function: $r(s_t, a_t) \ge 0$

(2)





0 Method overview



Encoder:

$$\boldsymbol{e}_{\phi}\left(\boldsymbol{z}_{t} \mid \boldsymbol{s}_{t}\right) \tag{3}$$

Aalto University School of Engineering

presented by Aidan Scannell

5/21 23rd September 2022

0 Method overview



Encoder:

$$\boldsymbol{e}_{\phi}\left(\boldsymbol{z}_{t} \mid \boldsymbol{s}_{t}\right) \tag{3}$$

Latent space dynamics:

$$m_{\phi}\left(z_{t} \mid z_{t-1}, a_{t-1}\right) \quad (4)$$

Aalto University School of Engineering

0 Method overview



Aalto University School of Engineering

• Optimality variable $\mathcal{O}_t \in \{0, 1\}, (1 = \text{optimal})^a$



^a[3] Toussaint. "Robot Trajectory Optimization Using Approximate Inference". 2009.



• Optimality variable $\mathcal{O}_t \in \{0, 1\}$, $(1 = optimal)^a$

Likelihood:

$$p(\mathcal{O}_{0:\infty} = 1 \mid \tau) \propto R(\tau) \tag{6}$$



where
$$\tau \triangleq (s_0, a_0, z_0, s_1, a_1, z_1, \cdots)$$

^a[3] Toussaint. "Robot Trajectory Optimization Using Approximate Inference". 2009.





■ Optimality variable O_t ∈ {0, 1}, (1 = optimal)^a
 ■ Likelihood:

$$p(\mathcal{O}_{0:\infty} = 1 \mid \tau) \propto R(\tau)$$
 (6)

where $\tau \triangleq (s_0, a_0, z_0, s_1, a_1, z_1, \cdots)$ Prior:

$$p_{\phi}(\tau) \triangleq p_0\left(s_0\right) \prod_{t=0}^{\infty} \underbrace{p\left(s_{t+1} \mid s_t, a_t\right)}_{\text{environment}} \underbrace{\pi_{\phi}\left(a_t \mid z_t\right)}_{\text{policy}} \underbrace{e_{\phi}\left(z_t \mid s_t\right)}_{\text{encoder}}$$

^a[3] Toussaint. "Robot Trajectory Optimization Using Approximate Inference". 2009.



0 Distribution over trajectories

$$p_{\phi}(\tau) \triangleq p_0(s_0) \prod_{t=0}^{\infty} \underbrace{p(s_{t+1} \mid s_t, a_t)}_{\text{environment}} \underbrace{\pi_{\phi}(a_t \mid z_t)}_{\text{policy}} \underbrace{e_{\phi}(z_t \mid s_t)}_{\text{encoder}}$$
(7)

Drawing samples from $p_{\phi}(\tau)$ requires interacting with the environment...



0 Distribution over trajectories

$$p_{\phi}(\tau) \triangleq p_0(s_0) \prod_{t=0}^{\infty} \underbrace{p(s_{t+1} \mid s_t, a_t)}_{\text{environment}} \underbrace{\pi_{\phi}(a_t \mid z_t)}_{\text{policy}} \underbrace{e_{\phi}(z_t \mid s_t)}_{\text{encoder}}$$
(7)

Drawing samples from $p_{\phi}(\tau)$ requires interacting with the environment...

Estimate $\mathbb{E}_{p(\tau)}[R(\tau)]$ with trajectories sampled from another distribution $q(\tau)$



0 Distribution over trajectories

$$p_{\phi}(\tau) \triangleq p_0(s_0) \prod_{t=0}^{\infty} \underbrace{p(s_{t+1} \mid s_t, a_t)}_{\text{environment}} \underbrace{\pi_{\phi}(a_t \mid z_t)}_{\text{policy}} \underbrace{e_{\phi}(z_t \mid s_t)}_{\text{encoder}}$$
(7)

Drawing samples from $p_{\phi}(\tau)$ requires interacting with the environment...

- Estimate $\mathbb{E}_{p(\tau)}[R(\tau)]$ with trajectories sampled from another distribution $q(\tau)$
- Lower bound expected return using $q(\tau)$

$$\log \mathbb{E}_{\rho(\tau)}[R(\tau)] \ge \mathbb{E}_{q(\tau)}[\log R(\tau) + \log p(\tau) - \log q(\tau)]$$
(8)



0 Variational distribution



Latent representations z_t are independent of states s_t for t > 0



0 Variational distribution



- Latent representations z_t are independent of states s_t for t > 0
- So can estimate bound for any policy using samples from latent-space model



0 Variational distribution



- Latent representations z_t are independent of states s_t for t > 0
- So can estimate bound for any policy using samples from latent-space model
 - Does not need access to high dimensional states!



$$\mathcal{L}_{\phi}^{K} \triangleq \mathbb{E}_{q_{\phi}^{K}(\tau)} \left[\left(\sum_{t=0}^{K-1} \gamma^{t} \tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) \right) + \gamma^{K} \log Q\left(s_{K}, a_{K}\right) \right]$$
where $\tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) = \underbrace{(1-\gamma)\log r\left(s_{t}, a_{t}\right)}_{\text{extrinsic term}} + \underbrace{\log e_{\phi}\left(z_{t+1} \mid s_{t+1}\right) - \log m_{\phi}\left(z_{t+1} \mid z_{t}, a_{t}\right)}_{\text{intrinsic}}$

Bound holds for any:



$$\mathcal{L}_{\phi}^{K} \triangleq \mathbb{E}_{q_{\phi}^{K}(\tau)} \left[\left(\sum_{t=0}^{K-1} \gamma^{t} \tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) \right) + \gamma^{K} \log Q\left(s_{K}, a_{K}\right) \right]$$

$$\text{where} \quad \tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) = \underbrace{\left(1 - \gamma\right) \log r\left(s_{t}, a_{t}\right)}_{\text{extrinsic term}} + \underbrace{\log e_{\phi}\left(z_{t+1} \mid s_{t+1}\right) - \log m_{\phi}\left(z_{t+1} \mid z_{t}, a_{t}\right)}_{\text{intrinsic}}$$

- Bound holds for any:
 - representation $e_{\phi}(z_t \mid s_t)$



$$\mathcal{L}_{\phi}^{K} \triangleq \mathbb{E}_{q_{\phi}^{K}(\tau)} \left[\left(\sum_{t=0}^{K-1} \gamma^{t} \tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) \right) + \gamma^{K} \log Q\left(s_{K}, a_{K}\right) \right]$$

$$\text{where} \quad \tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) = \underbrace{(1-\gamma) \log r\left(s_{t}, a_{t}\right)}_{\text{extrinsic term}} + \underbrace{\log e_{\phi}\left(z_{t+1} \mid s_{t+1}\right) - \log m_{\phi}\left(z_{t+1} \mid z_{t}, a_{t}\right)}_{\text{intrinsic}}$$

- Bound holds for any:
 - representation $e_{\phi}(z_t \mid s_t)$
 - latent-space model $m_{\phi}(z_{t+1} \mid z_t, a_t)$

$$\mathcal{L}_{\phi}^{K} \triangleq \mathbb{E}_{q_{\phi}^{K}(\tau)} \left[\left(\sum_{t=0}^{K-1} \gamma^{t} \tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) \right) + \gamma^{K} \log Q\left(s_{K}, a_{K}\right) \right]$$

$$\text{where} \quad \tilde{r}\left(s_{t}, a_{t}, s_{t+1}\right) = \underbrace{(1-\gamma) \log r\left(s_{t}, a_{t}\right)}_{\text{extrinsic term}} + \underbrace{\log e_{\phi}\left(z_{t+1} \mid s_{t+1}\right) - \log m_{\phi}\left(z_{t+1} \mid z_{t}, a_{t}\right)}_{\text{intrinsic}}$$

- Bound holds for any:
 - representation $e_{\phi}(z_t \mid s_t)$
 - latent-space model $m_{\phi}(z_{t+1} \mid z_t, a_t)$
 - policy $\pi_{\phi}(a_t \mid z_t)$

0 Practical algorithm - Aligned Latent Models

Algorithm 1 The ALM objective can be optimized with any RL algorithm. We present an implementation based on DDPG Lillicrap et al. (2015).

1: Initialize the encoder $e_{\phi}(z_t \mid s_t)$, model $m_{\phi}(z_{t+1} \mid z_t, a_t)$, policy $\pi_{\phi}(a_t \mid z_t)$, classifier $C_{\theta}(z_{t+1}, a_t, z_t)$, reward $r_{\theta}(z_t, a_t)$, Q-function $Q_{\theta}(z_t, a_t)$, replay buffer \mathcal{B}

2: for
$$n = 1, \cdots, N$$
 do do

- 3: Select action $a_n = \pi_{\phi}(a_n \mid e_{\phi}(s_n)) + \mathcal{N}$ using the current policy and exploration noise \mathcal{N} .
- 4: Execute action a_n and observe reward r_n and next state s_{n+1} .
- 5: Store transition (s_n, a_n, r_n, s_{n+1}) in \mathcal{B} ; sample length-K sequence $(s_i, a_i, s_{i+1})_{i=t}^{t+K-1} \sim \mathcal{B}$
- 6: Compute lower bound using off-policy actions: $\mathcal{L}_{e_{\phi},m_{\phi}}^{K}((s_{i},a_{i},s_{i+1})_{i=t}^{t+K-1}) \triangleright 7$
- 7: Update encoder and model by gradient ascent on off-policy lower bound: $\mathcal{L}_{e_{\phi},m_{\phi}}^{K}$
- 8: Compute lower bound using on-policy actions: $\mathcal{L}_{\pi_{\phi}}^{K}((s_{t}))$ \triangleright 8
- 9: Update policy by gradient ascent on on-policy lower bound: $\mathcal{L}_{\pi_{\phi}}^{K}$
- 10: Update classifier, Q-function and reward by gradient descent on: $\mathcal{L}_{C_{\theta}}, \mathcal{L}_{Q_{\theta}}, \mathcal{L}_{r_{\theta}} > 11, 9, 10$

0 Practical algorithm - Latent-space learning phase

$$\mathcal{L}_{\boldsymbol{e}_{\phi},\boldsymbol{m}_{\phi}}^{K}\left(\{\boldsymbol{s}_{i},\boldsymbol{a}_{i},\boldsymbol{s}_{i+1}\}_{i=t}^{t+K-1}\right) = \mathbb{E}_{\underbrace{\boldsymbol{e}_{\phi}\left(\boldsymbol{z}_{i=t}\mid\boldsymbol{s}_{t}\right), \underbrace{\boldsymbol{m}_{\phi}(\boldsymbol{z}_{i>t}\mid\boldsymbol{z}_{t:i-1},\boldsymbol{a}_{i-1})}_{\text{latent-space dynamics}}}\left[\gamma^{K}\underbrace{\boldsymbol{Q}_{\theta}\left(\boldsymbol{z}_{K},\boldsymbol{\pi}\left(\boldsymbol{z}_{K}\right)\right)}_{\textit{learned,latent}} + \sum_{i=t}^{t+K-1}\gamma^{i}\left(\underbrace{\boldsymbol{r}_{\theta}\left(\boldsymbol{z}_{i},\boldsymbol{a}_{i}\right)}_{\textit{learned,latent}} - \underbrace{\mathsf{KL}\left(\boldsymbol{e}_{\phi_{\text{targ}}}\left(\boldsymbol{z}_{i+1}\mid\boldsymbol{s}_{i+1}\right) \|\boldsymbol{m}_{\phi}\left(\boldsymbol{z}_{i+1}\mid\boldsymbol{z}_{i},\boldsymbol{a}_{i}\right)\right)}_{\text{latent-space consistency}}\right)\right]$$



0 Practical algorithm - Planning phase

$$= \text{Remeber the original objective } \mathcal{L}_{\phi}^{K} \triangleq \mathbb{E}_{q_{\phi}^{K}(\tau)} \left[\left(\sum_{t=0}^{K-1} \gamma^{t} \tilde{r} \left(s_{t}, a_{t}, s_{t+1} \right) \right) + \gamma^{K} \log Q\left(s_{K}, a_{K} \right) \right]$$
where $\tilde{r}\left(s_{t}, a_{t}, s_{t+1} \right) = \underbrace{(1-\gamma) \log r\left(s_{t}, a_{t} \right)}_{\text{extrinsic term}} + \underbrace{\log e_{\phi}\left(z_{t+1} \mid s_{t+1} \right) - \log m_{\phi}\left(z_{t+1} \mid z_{t}, a_{t} \right)}_{\text{intrinsic}}$

0 Practical algorithm - Planning phase

$$= \text{Remeber the original objective } \mathcal{L}_{\phi}^{K} \triangleq \mathbb{E}_{q_{\phi}^{K}(\tau)} \left[\left(\sum_{t=0}^{K-1} \gamma^{t} \tilde{r}(s_{t}, a_{t}, s_{t+1}) \right) + \gamma^{K} \log Q(s_{K}, a_{K}) \right]$$
where $\tilde{r}(s_{t}, a_{t}, s_{t+1}) = \underbrace{(1-\gamma) \log r(s_{t}, a_{t})}_{\text{extrinsic term}} + \underbrace{\log e_{\phi}(z_{t+1} \mid s_{t+1}) - \log m_{\phi}(z_{t+1} \mid z_{t}, a_{t})}_{\text{intrinsic}}$

Do not have access to s_{t+1} for on-policy actions, so can't calculate $e_{\phi}(z_{t+1} | s_{t+1})$

0 Practical algorithm - Planning phase

a Remeber the original objective
$$\mathcal{L}_{\phi}^{K} \triangleq \mathbb{E}_{q_{\phi}^{K}(\tau)} \left[\left(\sum_{t=0}^{K-1} \gamma^{t} \tilde{r}(s_{t}, a_{t}, s_{t+1}) \right) + \gamma^{K} \log Q(s_{K}, a_{K}) \right]$$
where $\tilde{r}(s_{t}, a_{t}, s_{t+1}) = \underbrace{(1-\gamma) \log r(s_{t}, a_{t})}_{\text{extrinsic term}} + \underbrace{\log e_{\phi}(z_{t+1} \mid s_{t+1}) - \log m_{\phi}(z_{t+1} \mid z_{t}, a_{t})}_{\text{intrinsic}}$
b not have access to s_{t+1} for on-policy actions, so can't calculate $e_{\phi}(z_{t+1} \mid s_{t+1})$
b Learn classifier $C_{\theta}(z_{t+1}, a_{t}, z_{t})$ to differentiate z_{t+1} from $e_{\phi}(z_{t+1} \mid s_{t+1})$ and $m_{\phi}(z_{t+1} \mid s_{t}, a_{t})$
 $\mathcal{L}_{\pi_{\phi}}^{K}(s_{t}) = \mathbb{E}_{q_{\phi}(z_{t:K}, a_{t:K} \mid s_{t})} \left[\sum_{i=t}^{t+K-1} \gamma^{t} \left(r_{\theta}(z_{i}, a_{i}) + c \cdot \log \frac{C_{\theta}(z_{i+1}, a_{i}, z_{i})}{1 - C_{\theta}(z_{i+1}, a_{i}, z_{i})} \right) + \gamma^{K} Q_{\theta}(z_{K}, \pi(z_{K})) \right]$



0 Experiments - Is the ALM objective useful?

Table 1: On the model-based benchmark from Wang et al. (2019), ALM outperforms model-based and modelfree methods on $\frac{4}{5}$ tasks, often by a wide margin. We report mean and std. dev. across 5 random seeds. T-Humanoid-v2 and T-Ant-v2 refer to the respective truncated environments.

	T-Humanoid-v2	T-Ant-v2	HalfCheetah-v2	Walker2d-v2	Hopper-v2
ALM(3)	5306 ± 437	4887 ± 1027	10789 ± 366	3006 ± 1183	2546 ± 1074
SAC-SVG(2) (Amos et al., 2020) SAC-SVG(3)	$\begin{array}{c} 501 \pm {\scriptstyle 37} \\ 472 \pm {\scriptstyle 85} \end{array}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c} 8752 \pm {\scriptstyle 1785} \\ 9220 \pm {\scriptstyle 1431} \end{array}$	$\begin{array}{c} 448 \pm {\scriptstyle 1139} \\ 878 \pm {\scriptstyle 1533} \end{array}$	$\begin{array}{ }\textbf{2852} \pm \textbf{361} \\ \textbf{2024} \pm \textbf{1981} \end{array}$
SVG(1) (Heess et al., 2015)	$811.8 \pm {\scriptstyle 241.5}$	$185 \pm {\scriptstyle 141}$	$336 \pm {\scriptscriptstyle 387}$	$ 252 \pm 48$	435 ± 163
SLBO (Luo et al., 2018)	1377 ± 150	200 ±40	$1097 \pm {\scriptstyle 166}$	$ $ 207 \pm 108	$805 \pm {\scriptstyle 142}$
TD3 (Fujimoto et al., 2018) SAC (Haarnoja et al., 2018)	$\begin{array}{c c} 147 \pm \textit{0.7} \\ 1470 \pm \textit{794} \end{array}$	$\begin{array}{c c} 870 \pm {\scriptstyle 283} \\ 548 \pm {\scriptstyle 146} \end{array}$	$\begin{array}{c} 3016 \pm {\scriptstyle 969} \\ 3460 \pm {\scriptstyle 1326} \end{array}$	$\begin{array}{c} \textbf{-516} \pm \textbf{812} \\ \textbf{166} \pm \textbf{1318} \end{array}$	$\begin{array}{c}1817 \pm 994 \\788 \pm 977\end{array}$

Policy return after 2e5 environment steps



0 Experiments - Can ALM achieve good performance without ensembles



(a) Sample efficiency comparison

(b) Wall-clock comparison

Figure 3: Good performance without ensembles. Our method (ALM) can (*Left*) match the sample complexity of ensembling-based methods (MBPO, REDQ) while (*Right*) requiring less compute. ALM performs updates $\sim 10 \times$ faster than MBPO. See Appendix Fig. 7 for results on other environments.



0 Experiments - Why does ALM work - Anaylzing the Q values



Figure 5: Analyzing Q-values. See text for details.

0 Experiments - Why does ALM work - Ablation experiments





0 Experiments - Why does ALM work - Analysing the learned representations





One objective for learning representations, latent-space dynamics, and policy

¹[2] Janner et al. "When to Trust Your Model: Model-Based Policy Optimization". 2019.
²[1] Chen et al. "Randomized Ensembled Double Q-Learning: Learning Fast Without a Model". 2021.





One objective for learning representations, latent-space dynamics, and policy
 Similar sample efficient to MBPO¹ and REDQ²

¹[2] Janner et al. "When to Trust Your Model: Model-Based Policy Optimization". 2019.
²[1] Chen et al. "Randomized Ensembled Double Q-Learning: Learning Fast Without a Model". 2021.



0 Summary

- One objective for learning representations, latent-space dynamics, and policy
- Similar sample efficieny to MBPO¹ and REDQ²
 - No ensembles leads to ~6x less wall-clock time

¹[2] Janner et al. "When to Trust Your Model: Model-Based Policy Optimization". 2019.
²[1] Chen et al. "Randomized Ensembled Double Q-Learning: Learning Fast Without a Model". 2021.



0 References

- [1] Xinyue Chen, Che Wang, Zijian Zhou, and Keith Ross. "Randomized Ensembled Double Q-Learning: Learning Fast Without a Model". In: *International Conference on Learning Representations*. arXiv, Mar. 2021. arXiv: 2101.05982 [cs].
- [2] Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. "When to Trust Your Model: Model-Based Policy Optimization". In: *Advances in Neural Information Processing Systems*. Vol. 32. 2019.
- [3] Marc Toussaint. "Robot Trajectory Optimization Using Approximate Inference". In: *International Conference on Machine Learning*. 2009.



0 Difference between theory and experiments

They omit $\log r(s_t, a_t)$ from Equations 7 and 8



0 Difference between theory and experiments

- They omit $\log r(s_t, a_t)$ from Equations 7 and 8
 - because it decreased performance



0 Difference between theory and experiments

- They omit $\log r(s_t, a_t)$ from Equations 7 and 8
 - because it decreased performance
- They scale the KL term in policy objective (Equation 7)



More moving parts than model-free algorithms

My thoughts



- More moving parts than model-free algorithms
 - **Target network for** e_{ϕ}

My thoughts



More moving parts than model-free algorithms

- **Target network for** e_{ϕ}
- Target network for Q_{θ}



More moving parts than model-free algorithms

- **Target network for** e_{ϕ}
- Target network for Q_{θ}

Classifer C_{θ} for estimating $\log e_{\phi}(z_{t+1} \mid s_{t+1}) - \log m_{\phi}(z_{t+1} \mid s_t, a_t)$



More moving parts than model-free algorithms

- Target network for *e*_φ
- Target network for Q_{θ}
- Classifer C_{θ} for estimating $\log e_{\phi}(z_{t+1} \mid s_{t+1}) \log m_{\phi}(z_{t+1} \mid s_t, a_t)$
- Gap between theoretical analysis and practical implementation (function approximation etc)



More moving parts than model-free algorithms

- Target network for *e*_φ
- Target network for Q_{θ}
- Classifer C_{θ} for estimating $\log e_{\phi}(z_{t+1} \mid s_{t+1}) \log m_{\phi}(z_{t+1} \mid s_t, a_t)$
- Gap between theoretical analysis and practical implementation (function approximation etc)

My thoughts

■ Fig 4 - Still good performance when removing the latent-space dynamics????



More moving parts than model-free algorithms

- Target network for *e*_φ
- Target network for Q_{θ}
- Classifer C_{θ} for estimating $\log e_{\phi}(z_{t+1} \mid s_{t+1}) \log m_{\phi}(z_{t+1} \mid s_t, a_t)$
- Gap between theoretical analysis and practical implementation (function approximation etc)

- Fig 4 Still good performance when removing the latent-space dynamics????
 - I'm a little confused about this as learning C_{θ} in Eq 7 requires access to m_{θ} ?



More moving parts than model-free algorithms

- Target network for *e*_φ
- Target network for Q_{θ}
- Classifer C_{θ} for estimating $\log e_{\phi}(z_{t+1} \mid s_{t+1}) \log m_{\phi}(z_{t+1} \mid s_t, a_t)$
- Gap between theoretical analysis and practical implementation (function approximation etc)

- Fig 4 Still good performance when removing the latent-space dynamics???
 I'm a little confused about this as learning C_θ in Eq 7 requires access to m_θ?
- This method encodes exploration via expectation over latent trajectories $\{z_t\}$?



More moving parts than model-free algorithms

- Target network for *e*_φ
- Target network for Q_{θ}
- Classifer C_{θ} for estimating $\log e_{\phi}(z_{t+1} \mid s_{t+1}) \log m_{\phi}(z_{t+1} \mid s_t, a_t)$
- Gap between theoretical analysis and practical implementation (function approximation etc)

- Fig 4 Still good performance when removing the latent-space dynamics???
 I'm a little confused about this as learning C_θ in Eq 7 requires access to m_θ?
- This method encodes exploration via expectation over latent trajectories $\{z_t\}$?
 - But the distribution over z_t does not represent epistemic uncertainty...

