

The idea is to compose a policy with following form:

$$\pi_{\theta}(a|s) \doteq \mathcal{N}(a | \text{Mean}_{\theta}(s), \text{DiagVar}_{\theta}) \quad (0)$$

$$\text{Where } \text{Mean}_{\theta}(s) = \alpha_{TV} \cdot a^{\text{ref}}(s) + (1 - \alpha_{TV}) \cdot a_{\theta}^{\text{explore}}(s) \quad (1)$$

- $a^{\text{ref}}(s)$ comes directly from a planning method on dynamics model, while $a_{\theta}^{\text{explore}}(s)$ is a neural network policy to be aggregated to the reference policy $a^{\text{ref}}(s)$ with a weight of $1 - \alpha_{TV}$.
- α_{TV} denotes a measure of Trust Value of reference policy.
- Trust Value should be dependent on the discrepancy between actual reward in real environment $r(s, a)$ and simulated reward from dynamics model $\hat{r}(s, a)$. So, α_{TV} should depend on $r^{\text{diff}} = |r(s, a) - \hat{r}(s, a)|$, which can be trained online in a supervised way.

- Supervised learning for r_ϕ^{diff} with parameter ϕ
- Then, α_{TV} , for example, could be formulated as:

$$\alpha_{TV}(s, a) = \exp(-r_\phi^{\text{diff}}(s, a)) \quad (2)$$

- $r_{\text{label}}^{\text{diff}}$ could be defined as normalized:

$$r_{\text{label}}^{\text{diff}}(s, a) = \frac{|r(s, a) - \hat{r}(s, a)|}{|r(s, a)|} \quad (3)$$

- α_{TV} could also be other functions of $r_\phi^{\text{diff}}(s, a)$.
- In summary, total loss could be:

$$\min_{\theta, \phi} \mathcal{L}(\theta, \phi) = -\mathcal{L}(\theta) + \beta \cdot \mathcal{L}(\phi), \text{ where } \beta \text{ is a hyperparameter.}$$

$$\begin{cases} \mathcal{L}(\theta) = \mathbb{E}_{s \sim \rho_{\text{old}}, a \sim \pi_{\theta_{\text{old}}}} \left[\frac{\pi_{\theta_{\text{new}}}(\mathbf{a}|s)}{\pi_{\theta_{\text{old}}}(\mathbf{a}|s)} A_{\theta_{\text{old}}}(s, a) \right] & (4) \\ \mathcal{L}(\phi) = \| r_\phi^{\text{diff}}(s, a) - r_{\text{label}}^{\text{diff}}(s, a) \|^2 & (5) \end{cases}$$

Discussions of $r_{\phi}^{\text{diff}}(s,a)$

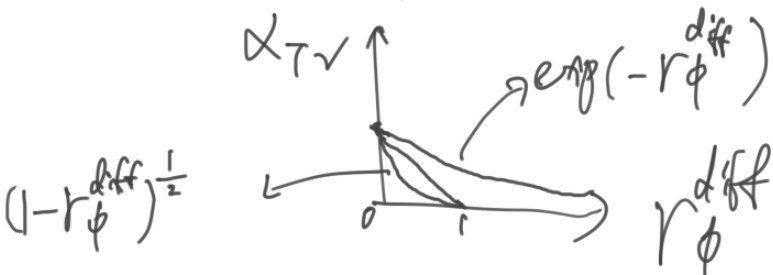
- ϕ Network learns $|r(s,a) - \hat{r}(s,a)|$?
or by normalizing $\frac{|r(s,a) - \hat{r}(s,a)|}{|\hat{r}(s,a)|}$?
if normalize, what is denominator? $|\hat{r}(s,a)|$ or $|r(s,a)|$ or?

- if $r(s,a) > \hat{r}(s,a)$, should we assume $\alpha_{TV}(r_{\phi}^{\text{diff}}) = 1$?
- Should we consider $r_{\phi}^{\text{diff}} \in [0, +\infty)$, or $\in [0, 1]$, or $\in (-\infty, +\infty)$
- Exact function α_{TV} of r_{ϕ}^{diff} ? $\alpha_{TV} = \exp(-r_{\phi}^{\text{diff}})$?

or other functions such as

$$\alpha_{TV} = 1 - r_{\phi}^{\text{diff}} ?$$

$$\alpha_{TV} = (1 - r_{\phi}^{\text{diff}})^p$$



Discussions of $\alpha_{TV} \cdot a^{ref}$

- more generally, linear transformation of a^{ref} ? i.e. $\alpha_{TV} \cdot \underline{T} \cdot a^{ref}$
where T is a transformation matrix.

- Most generally, $\alpha_{TV} \cdot f_w(a^{ref})$, where f_w is a NN with w as parameter.

