

CIM-PPO: Proximal Policy Optimization with Liu-Correntropy Induced Metric[†]

Yunxiao Guo, Han Long*, Xiaojun Duan*, Kaiyuan Feng, Maochu Li, Xiaying Ma

Abstract—As an algorithm based on deep reinforcement learning (DRL), Proximal Policy Optimization (PPO) performs well in many complex tasks. According to the mechanism of penalty in a surrogate objective, PPO can be divided into PPO with KL divergence (KL-PPO) and PPO with Clip function (Clip-PPO). In this paper, the authors analyze the effect of asymmetry of KL divergence on PPO's objective function and give an inequality that can indicate when the asymmetry will affect the efficiency of KL-PPO first. Next, by replacing KL divergence with Correntropy Induced Metric (CIM) in the objective function of PPO, the authors construct the CIM-PPO algorithm. The CIM-PPO algorithm is more robust because the correntropy induced entropy has robust performance in the distance measurement of random variables for signal processing problems under non-Gaussian noise conditions. Finally, the authors design experiments based on OpenAI gym to test the effectiveness of the CIM-PPO algorithm. It shows that CIM-PPO performs better than KL-PPO and Clip-PPO.

Index Terms—KL divergence, Proximal policy optimization (PPO), Correntropy induced metric (CIM), Surrogate objective, Deep reinforcement learning (DRL)

I. INTRODUCTION

AS a significant part of machine learning, reinforcement learning (RL) aims to obtain the optimal policy of the agents interacting with the environment. Traditional RL algorithms face lots of problems, such as the curse of dimensionality, continuous control, and large-scale calculation. In recent years, many reinforcement learning algorithms have been proposed to solve the problems of traditional reinforcement learning [1]–[4].

The most famous method to solve the curse of dimensionality is deep reinforcement learning (DRL). DRL combines the advantages of deep learning and reinforcement learning [5], [6]. The DRL effectively solved the curse of dimensionality by applying deep neural networks and it handles well when the input state space is an image [5]–[8]. In recent years, DRL has been widely used in various fields, such as video games [9], visual object tracking [10], scheduling optimization [11], autonomous driving [12], natural language processing [13].

In solving continuous control problems, policy gradient is a common method. It tries to reach the optimal policy by following the future reward function's gradient. An obvious

benefit of policy gradient methods is that they can ensure improved performance at every update step in theoretic. However, policy gradient is difficult to do policy update and policy evaluation at the same time. Therefore, the Actor-Critic struct is introduced in the policy gradient algorithms to solve this dilemma. By using two neural networks: actor-network and critic network, the policy updates, and policy evaluation can be done independently [14], [15].

Actor-critic struct methods also meet some problems in applications. For one thing, it can not guarantee the future reward function is monotone increasing after an update. For another, it is hard to converge on locally optimal policy while running into a plateau area. The similar question is discussed in [16] and [17]. In [17], the updated policy's performance was defined as the old policy's performance plus a nonnegative increasement. [16] uses the natural gradient instead of the standard gradient in the policy gradient part. **Trust Region Policy Optimization (TRPO)** was put forward under their inspiration. To ensure the increase of future reward function and the robust policy update step in the future, it adds a trust region constraint (KL divergence) instead of penalty to optimize the surrogate objective. Meanwhile, for a similar purpose, Levine and Abbeel [18] used the KL divergence constraint and dynamics model to restrict policy staying at the valid region, which can guarantee the increase of reward function. This method has achieved improvement, especially in robotic locomotion experiments by using the MuJoCo simulator [19].

Although the importance sampling and approximate techniques were used, the methods based on trust region optimization are not always perfect. There is still a crop of problems. One of the questions is the large-scale learning tasks. By using the conjugate gradient method in the practical algorithm, a large quantity of Fisher-vectors' multiplication will be done. The efficiency of TRPO will be descent. One way to solve this problem is to use the clip function, and this method was firstly used in **Proximal Policy Optimization (PPO with clip)** [20]. In PPO with clip, the KL divergence is replaced by clip function, which is much easier to calculate. Besides, the constraint problem is also transformed into an unconstrained optimization problem. Thus, clip-PPO can deal with large-scale tasks with fewer calculations, and the results perform well. There is also another PPO algorithm which is mentioned in the article [20]: PPO with Adaptive KL Penalty Coefficient (KL-PPO), but the performance is worse than PPO-clip.

In recent years, many researchers have analyzed the clip-PPO algorithm from different perspectives and improved it. By utilizing the policy information in the process of the

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

Han Long and Xiaojun Duan are the corresponding authors. Yunxiao Guo et. al is with the College of Liberal Arts and Sciences, National University of Defense Technology, Changsha, 410073, Hunan, China. email: longhan@nudt.edu.cn

value function update, PPO with policy feedback [2] put forward a new mechanism, which is more efficient than clip-PPO. Proximal Policy Optimization with Relative Pearson Divergence [21] proposed a PPO algorithm based on relative Pearson divergence, and it performed as well as or better than clip-PPO. Yuhu Cheng et.al proposed a novel PPO algorithm with the first-order policy gradient, which is called authentic boundary PPO [22]. It effectively improves the learning stability and accelerates the convergence speed. However, these improvements are based on the clip-PPO algorithm, and there are few improvements to KL-PPO. In the process of one-step update, KL-PPO requires more computation than clip-PPO, so the convergence time of KL-PPO tends to be longer than clip-PPO, and is less stable than clip-PPO. Nevertheless, KL-PPO has a deeper theory than Clip-PPO, it uses KL divergence to measure the degree of difference between the new policy and the old policy, and adds it to the surrogate objective as the penalty. Then, the adaptive method is used to estimate the penalty coefficient. However, the KL-PPO, as a more theoretical and interpretable algorithm, has received little attention from the current research. Therefore, this paper analyzes its convergence rate and gives an inequality to prove that when the inequality is satisfied, the KL-PPO efficiency will decrease. Finally, we proposed a more stable and efficient algorithm, and experiments have been done to verify its properties. The main contributions of this paper are shown as follows:

- 1) We analyzed the properties of KL divergence and its influence on KL-PPO, and found that its asymmetry, large amount of computation, and difficulty to calculate in specific cases would restrict the performance of the KL-PPO algorithm. Specifically, an inequality is given through theoretical derivation. When the inequality is true, the update efficiency and robustness of the KL-PPO algorithm will be reduced due to its penalty term.
- 2) We proposed the CIM-PPO algorithm, which canceled the original KL divergence adaptive adjustment coefficient mechanism of KL-PPO. By adopting the Correntropy mechanism, the correntropy induced entropy(CIM) is used as a new penalty term for the updating of constraint policies. In the distance measurement of random variables under non-Gaussian noise conditions, the CIM has robust performance.

This paper is divided into six sections. The first section is a brief overview of the policy optimization methods and their development. The second section illustrates the foundation of reinforcement learning. The third section demonstrates the development of the PPO algorithm and gives the theorems to show the shortages of KL-PPO. The fourth section proposes a new estimate method to improve the policy optimization methods. The fifth section analyses the consequences of experiments. Our discussions are brought in the final section.

II. PRELIMINARIES

A reinforcement learning task that satisfies the Markov property is called Markov Decision Process (MDP) [23]–[25]. Especially, if the state and action spaces are both finite, it

calls finite Markov Decision Process (finite MDP). If a finite MDP that expected reward is consistent of infinite steps and with a discount factor, it is called infinite-horizon discounted Markov Decision Process(infinite-horizon MDP). A MDP is defined by the tuple $(\mathcal{S}, \mathcal{A}, P, r, d_0, \gamma)$ where \mathcal{S} is a finite set of states, \mathcal{A} is a finite set of actions, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the transition probability distribution, $r : \mathcal{S} \rightarrow \mathbb{R}$ is the reward function, $d_0 : \mathcal{S} \rightarrow \mathbb{R}$ is the initial state distribution, and $\gamma \in (0, 1)$ is the discount factor.

The policy $\pi(a_t|s_t)$ is the distribution of action a_t under the state s_t . The target of a continuous reinforcement learning task is to make the discounted return:

$$R_t = \sum_{t=0}^{\infty} \gamma^t r(s_t) \quad (1)$$

when it reaches maximum, we take the expectation of the discounted return under the policy π and w.r.t $\eta(\pi)$:

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

where $s_0 \sim d_0(s_0)$, $a_t \sim \pi(a_t|s_t)$, $s_{t+1} \sim P(s_{t+1}|s_t, a_t)$

To estimate the potential return of state s_t , define value function $V_\pi(s_t)$ as:

$$V_\pi(s_t) = \mathbb{E}_{a_t, s_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right] \quad (2)$$

Similarly, define the state-action value function $Q_\pi(s_t, a_t)$, which represents the expected return that the agent starting from s_t , taking the action a_t , and thereafter following policy π :

$$Q_\pi(s_t, a_t) = \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right], \quad (3)$$

And then, define the difference between value function and state-action function as the advantage function $A_\pi(s_t, a_t)$, w.r.t:

$$A_\pi(s_t, a_t) = Q_\pi(s_t, a_t) - V_\pi(s_t) \quad (4)$$

Expand $Q_\pi(s_t, a_t)$:

$$\begin{aligned} Q_\pi(s_t, a_t) &= \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} \left[r(s_t) + \gamma \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l+1}) \right] \middle| s_t, a_t \right] \\ &= \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} [r(s_t) + \gamma V_\pi(s_{t+1}) | s_t, a_t] \end{aligned} \quad (5)$$

Then the relationship between $V_\pi(s_t)$ and $Q_\pi(s_t, a_t)$ is as follows:

$$Q_\pi(s_t, a_t) = \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} [r(s_t)] + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} [V_\pi(s_{t+1})] \quad (6)$$

We call the states-actions series: $\{s_0, a_1, s_1, \dots\}$ as trajectory, w.r.t $\tau(s_0)$. Particularly, if a trajectory is sampled by a policy π , we write it as $\tau|\pi$. Then, utilize equation (5), $A_\pi(s_t, a_t)$ can be rewritten as:

$$A_\pi(s_t, a_t) = \mathbb{E}_{\tau(s_{t+1})} [\gamma V_\pi(s_{t+1}) + r(s_{t+1}) - V_\pi(s_t)] \quad (7)$$

Moreover, use $A_\pi(s_t, a_t)$ substitute $r(s_t)$ in equation (1) and take expectation, we have:

$$\begin{aligned} & \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k A_\pi(s_{t+k}, a_{t+k}) \right] \\ &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k (r(s_{t+k}) + \gamma V_\pi(s_{t+1+k}) - V_\pi(s_{t+k})) \right] \\ &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r(s_{t+k}) + \sum_{k=0}^{\infty} \gamma^{t+1+k} V_\pi(s_{t+1+k}) - \sum_{k=0}^{\infty} \gamma^k V_\pi(s_{t+k}) \right] \\ &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r(s_{t+k}) - V_\pi(s_0) \right] \\ &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r(s_{t+k}) \right] - \mathbb{E} [V_\pi(s_0)] \end{aligned} \quad (8)$$

Noticed, $\mathbb{E} [V_\pi(s_0)] = \mathbb{E}_{s_0, a_0, \dots} [\sum_{t=0}^{\infty} \gamma^t r(s_t)] = \eta(\pi)$. If we generate trajectory from another policy π^* by using $A_\pi(s_t, a_t)$ in calculate advantage function, we have:

$$\begin{aligned} & \mathbb{E}_{\tau|\pi^*} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] \\ &= \mathbb{E}_{\tau|\pi^*} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right] - \mathbb{E}_{s_0} [V_\pi(s_0)] \\ &= \eta(\pi^*) - \eta(\pi) \end{aligned} \quad (9)$$

Then, we have:

$$\eta(\pi^*) = \eta(\pi) + \mathbb{E}_{\tau|\pi^*} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] \quad (10)$$

To study the properties of the latter term, define the γ -discounted future state distribution as $d_{\pi^*, \mathcal{S}}(s)$:

$$d_{\pi^*, \mathcal{S}}(s) = \sum_{t=0}^{\infty} \gamma^t P(s_t = s) \quad (11)$$

$d_{\pi^*, \mathcal{S}}(s)$ can also reflect the visitation frequencies of state s under policy π^* . So, considering equation (10) and equation (11), a more deterministic result can be deducted:

$$\eta(\pi^*) = \eta(\pi) + \sum_s d_{\pi^*, \mathcal{S}}(s) \sum_a \pi^*(a|s) A_\pi(s, a). \quad (12)$$

The optimal policy π^* satisfies $\eta(\pi^*) > \eta(\pi)$ (Optimality), the latter term in above equation satisfying nonnegative condition [25]: $\sum_s d_{\pi^*, \mathcal{S}}(s) \sum_a \pi^*(a|s) A_\pi(s, a) \geq 0$. Moreover, if we relax the requirement, take the policy better than π : $\tilde{\pi}$, we will gain the same nonnegative result. Considering the context, if we sample actions from π^* , but sample states from $d_{\pi, \mathcal{S}}$, the right hand of equation (12) will be easy to optimize [26]. So, we define the local approximation function as:

$$\mathcal{L}_\pi(\pi^*) = \eta(\pi) + \sum_s d_{\pi, \mathcal{S}}(s) \sum_a \pi^*(a|s) A_\pi(s, a) \quad (13)$$

In another view, \mathcal{L}_π using the previous state distribution $d_{\pi, \mathcal{S}}$ means we ignoring the change of state distribution in the policy updating process.

III. POLICY OPTIMIZATION

A. Conservative Policy Iteration

Before neural network methods are considered, there are two standard methods to maximize the reward of tasks: greedy dynamic programming and policy gradient methods. Kakade and Langford [17] regard these traditional methods can not always perform well in the three following questions:

- 1) Is there a way to measure performance that can ensure the policy is improved with each update.
- 2) Can the algorithm identify the particular update that improves the measure mentioned above.
- 3) How does the algorithm evaluate the performance after a reasonable number of updates.

So, they put forward **Conservative Policy Iteration(CPI)** which can answer the three above questions perfectly.

CPI aims to find an "approximate" optimal policy, which can be generated by an approximate greedy policy chooser. Meanwhile, for the generated policy π' and a constant ϵ , the greedy policy chooser can guarantee the new policy advantage is larger than the old one with ϵ at least.

Moreover, using the mixture of the present policy π and a better policy π' : $\pi_{new} = (1 - \alpha)\pi + \alpha\pi'$ (where $\alpha \in [0, 1]$) as update rule and expanding the policy according to the law of total probability, Kakade and Langford derivative the follow Inequality:

$$\eta(\pi_{new}) - \eta(\pi) \geq \frac{\alpha}{1 - \gamma} (\mathbb{E}_{s \sim d_{\pi, \mathcal{S}}, a \sim \pi'} [A_\pi] - \frac{2\alpha\gamma\epsilon}{1 - \gamma(1 - \alpha)}) \quad (14)$$

It provides an ideal that finds a judicious α to make sure the right hand of the above inequality is nonnegative. Consequently, CPI can gain an improved policy in a short period. However, the improvement is limited, at the same time, they studied the limitation of the policy improvement, while the ϵ is too large, the improvement of the optimal policy $\tilde{\pi}$ will be constrained by the following inequality:

$$\eta(\tilde{\pi}) - \eta(\pi) \leq \frac{\epsilon}{(1 - \gamma)} \left\| \frac{d_{\tilde{\pi}, \mathcal{S}}}{d_{\pi, \mathcal{S}}} \right\|_\infty \quad (15)$$

Where $\|\cdot\|_\infty$ means the ℓ_∞ -norm.

This is a strategy optimization method that optimizes the surrogate function on one side of the inequality so that the function to be optimized on the other side can be enhanced. It has inspired many researchers.

B. Trust Region Policy Optimization

Inspired by CPI, by minimizing a certain surrogate objective function, J.Schulman et. al proposed an algorithm: TRPO [26], which can answer the three questions that are put forward in CPI as well as, and forced update step sizes in the trust region.

Combine equation (13) and equation (14)(or see [26] Appendix B, which proved it by perturbation theory), we can obtain the following inequality:

$$\eta(\pi_{new}) \geq \mathcal{L}_{\pi_{old}}^{CPI}(\pi_{new}) - \frac{2\epsilon\gamma}{(1-\gamma)^2} \alpha^2 \quad (16)$$

where $\epsilon = \max_s \mathbb{E}_{a \sim \pi'(a|s)} [A_\pi(s, a)]$.

Initially, TRPO considered the total variation divergence as $D_{TV}(\pi_{new} \parallel \pi_{old})$, which described the distance measure between the old policy and the new policy.

Then, analyze the relationship between the total variation divergence and the KL-divergence:

$$D_{TV}(p \parallel q)^2 \leq D_{KL}(p \parallel q) \quad (17)$$

$$D_{KL}^{\max}(\pi, \pi_{new}) = \max_s D_{KL}(\pi(\cdot|s) \parallel \pi_{new}(\cdot|s))$$

In objective 16, a slightly larger α can more likely ensure the improvement of the new policy. Therefore, the new inequality can be written as:

$$\eta(\pi_{new}) \geq \mathcal{L}_\pi(\pi_{new}) - \frac{4\epsilon\gamma}{(1-\gamma)^2} D_{KL}^{\max}(\pi, \pi_{new}) \quad (18)$$

Obviously, the target of RL is to maximize the expected return $\eta(\pi)$. According to the inequality 18, if the right part of the inequality gets increasing, the left part has had to ascend. Denote parameter $C = \frac{4\epsilon\gamma}{(1-\gamma)^2}$, our target can be written as the maximum problem:

$$\underset{\pi_{new}}{\text{maximize}} [\mathcal{L}_\pi(\pi_{new}) - CD_{KL}^{\max}(\pi, \pi_{new})] \quad (19)$$

In practice tasks, if directly optimal penalty coefficient C above the problem, the step sizes would be very small. In other words, the rate of policy convergence is going to be very slow. To accelerate the convergence, TRPO transforms Eq. (19) to a constraint problem:

$$\underset{\pi_{new}}{\text{maximize}} \mathcal{L}_\pi^{TRPO}(\pi_{new}) \quad (20)$$

subject to $D_{KL}^{\max} \rho_\pi(\pi, \pi_{new}) \leq \delta$.

Noticed that, the constraint $D_{KL}^{\max} \rho_\pi(\pi, \pi_{new})$ is hard to achieve because it will bring about large-scale calculation. Meanwhile, in order to take the information of samples, $D_{KL}^{\max} \rho_\pi(\pi, \pi_{new})$ is replaced by $\bar{D}_{KL}^{\rho_\pi}(\pi, \pi_{new})$, the average value of KL divergence. Now, if we utilize expectation's term to illustrate the problem, some statics numerical methods can be utilized. In TRPO, Monte Carlo methods and importance sampling are introduced to reduce calculation:

$$\underset{\pi_{new}}{\text{maximize}} \mathcal{L}_\pi^{TRPO}(\pi_{new}) = \mathbb{E}_{s \sim d_\pi, a \sim q} \left[\frac{\pi_{new}(a|s)}{\pi(a|s)} A_\pi(s, a) \right] \quad (21)$$

subject to $\mathbb{E}_{s \sim d_\pi} [D_{KL}(\pi(\cdot|s) \parallel \pi_{new}(\cdot|s))] \leq \delta$.

Although the static techniques are took, a large-scale calculation still remains, noticed the KL divergence calculate form:

$$D_{KL}(P \parallel Q) = \int p(x) \ln \frac{p(x)}{q(x)} dx \quad (22)$$

It means that we have calculated a KL-divergence in each update step. So the problem of heavy computation has not been resolved. Therefore, some numerical methods should be taken. Firstly, TRPO constructs a Fisher Information Matrix(FIM): A by analytically computing the Hessian of the KL-divergence. Then, using gradient in a linear approximation to \mathcal{L} and a quadratic approximation to the KL, and parameterizing the policy π to π_θ , we have the final optimal problem:

$$\underset{\theta_{new}}{\text{maximize}} \left[\nabla_\theta \mathcal{L}_\theta^{TRPO}(\theta_{new}) \Big|_{\theta_{new}=\theta} \cdot (\theta_{new} - \theta) \right] \quad (23)$$

subject to $\frac{1}{2}(\theta - \theta_{new})^T F(\theta)(\theta - \theta_{new}) \leq \delta$,

where $F(\theta)_{ij} = \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \mathbb{E}_{s \sim \rho_\pi} [D_{KL}(\pi(\cdot|s, \theta) \parallel \pi(\cdot|s, \theta_{new}))]$

By optimizing the above constraint problem, TRPO makes the DRL tasks faster and easier and increases performance, especially in continuous tasks. Besides, TRPO performance well in some wide variance tasks, like learning simulated robotic tasks and playing Atari games.

C. Proximal Policy Optimization

In TRPO, quantities of numerical techniques are took to make calculate easier, such as approximating by Fisher information matrix(FIM). In the article [26], authors even attempt to use the ℓ_2 norm: $\frac{1}{2} \|\theta - \theta_{new}\|_2^2 \leq \delta$ as the constraint.

In the methodology of solving 16, TRPO thinks it is hard to control the update step if a penalty for the coefficient is used in the objective. On the contrary to TRPO, PPO accepts the penalty in the objective. In the article [20], two algorithms are proposed: PPO with Adaptive KL Penalty Coefficient(PPO-KL) & PPO with clip function(PPO-clip)

Adaptive KL-PPO optimize the KL-penalty objective:

$$\mathcal{L}^{KL}(\pi) = \mathbb{E}_\pi \left[\frac{\pi_{new}(a|s)}{\pi(a|s)} A_\pi(s, a) - \beta D_{KL}(\pi(\cdot|s) \parallel \pi_{new}(\cdot|s)) \right] \quad (24)$$

Where $\mathbb{E}_\pi \left[\frac{\pi_{new}(\cdot)}{\pi(\cdot)} \right]$ represent importance sampling. The pseudocode of KL-PPO see in algorithm 1

Algorithm 1 PPO with KL divergence

Initialize π_0, d_{targ} .
for $i = 0, 1, 2, \dots$ until convergence **do**
 Optimize follow estimation by minibatch SGD:
 Estimate $\hat{\sigma}$ by sampling old policy and new policy
 $\pi_{i+1} = \arg \max_{\pi} \mathbb{E}_{\pi_i} \left[\frac{\pi(a|s)}{\pi_i(a|s)} A_{\pi_i} - \beta D_{\text{KL}}(\pi_i(\cdot|s) \parallel \pi(\cdot|s)) \right]$
 Compute $d = D_{\text{KL}}(\pi_i(\cdot|s) \parallel \pi(\cdot|s))$.
 if $d < d_{targ}/1.5$
 $\beta = \beta/2$
 if $d > d_{targ} \times 1.5$
 $\beta = \beta \times 2$
end for

Using the adaptive coefficient β , if KL-divergence gets bigger, β will be shrunk in the next iteration; while KL-divergence gets smaller, β will be raised in the next iteration as well. Therefore, the penalty will be tending to be stationary in the iteration, the update epoch steps will be controlled well. PPO-clip is more innovative, the penalty as follow:

$$\mathcal{L}^{CLIP}(\pi) = \mathbb{E}_{\pi} \left[\min \left(\frac{\pi_{new}(a|s)}{\pi(a|s)} A_{\pi}, \text{clip} \left(\frac{\pi_{new}(a|s)}{\pi(a|s)}, 1 - \varepsilon, 1 + \varepsilon \right) A_{\pi} \right) \right] \quad (25)$$

where a clip function is defined as:

$$\text{clip}(x, 1 - \varepsilon, 1 + \varepsilon) = \begin{cases} 1 + \varepsilon, & x \geq 1 + \varepsilon \\ 1 - \varepsilon, & x \leq 1 - \varepsilon \\ x, & \text{other} \end{cases} \quad (26)$$

Clip function can constraint an input x in the interval $[1 - \varepsilon, 1 + \varepsilon]$ with a negligible computation. With Clip, PPO performance is better compared to the past method, especially in large-scale RL tasks that Deep Learning is based on.

In PPO-clip, the KL- divergence in the surrogate objective is replaced by a clip function, although there are no complete theories to support this replacement and it reduces the interpretability. PPO-clip still can perform strikingly. Nonetheless, it inspires us, can some theoretical method be taken to make the algorithm perfect? Can we analyze the significance of KL-divergence in the whole model, and put forward a new algorithm?

In the next part, we will use the theoretical method to analyze the current algorithm and discuss the problems of KL-PPO.

D. The Problems with Current Method

The PPO or TRPO are the optimizations of surrogates, which maximize the surrogate objective \mathcal{L}_{π} to ensure the discount reward can be approximated the maximum. Compare to TRPO, PPO-KL transforms the constraint condition to a penalty function multiplying a penalty factor. However, the introduction of the penalty in PPO-KL will cause low efficient updates. To avoid the policy vibrating in each update, KL divergence was introduced as a constraint. Nevertheless, the

KL divergence is asymmetric and does not obey triangle inequality, in the case of normal distribution, the asymmetry of KL divergence will increase as the means and variances turn to be different, in the condition of two one dimensional normal distributions: p, q , the asymmetry can be revealed:

$$\begin{aligned} D_{\text{KL}}(p \parallel q) - D_{\text{KL}}(q \parallel p) &= \log \left(\frac{\sigma_2^2}{\sigma_1^2} \right)^2 + \frac{(\sigma_1^2 - \sigma_2^2)[(\mu_1 - \mu_2)^2 + (\sigma_1^2 + \sigma_2^2)]}{2\sigma_1^2\sigma_2^2} \end{aligned}$$

See appendix A for proof.

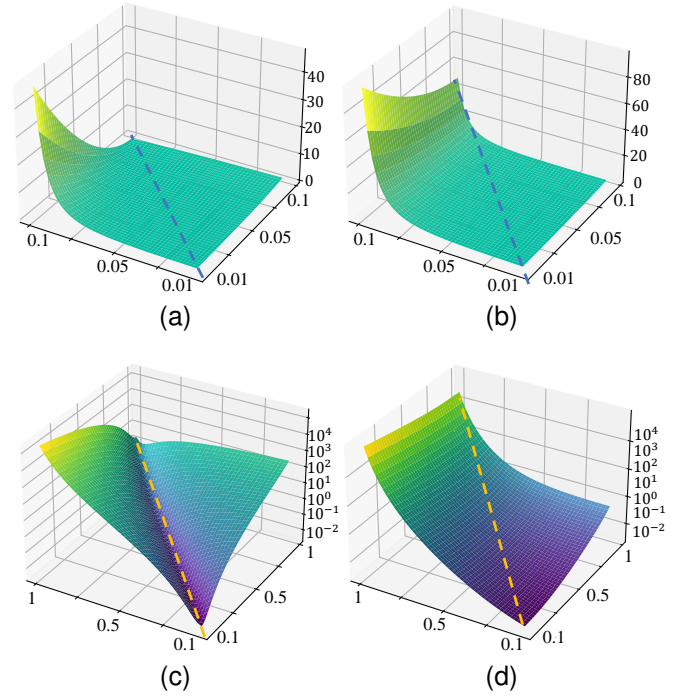


Fig. 1. The 3D-plots of KL divergence with different parameters (Different color represent different value of $D_{\text{KL}}(\pi_1 \parallel \pi_2)$): (a): means: $\mu_1 = \mu_2 = 1$ and variance variant from 0.01 to 0.1; (b): means: $\mu_1 = 1, \mu_2 = 1.1$ and variance variant from 0.01 to 0.1; (c): means: $\mu_1 = 1, \mu_2 = 1.1$ and variance variant from 0.01 to 1; (d): means: $\mu_1 = 1, \mu_2 = 2$ and variance variant from 0.01 to 1 and 1 to 10. The dotted lines just used to help our to see the asymetry more obviously. In the 3D-plots, we can obviously observed the color is not symmetric along the dotted diagonal.

In practical PPO and TRPO, the policy is usually parametrized by Deep Neural Networks (DNNs) [27]. For continuous action space tasks, it is standard to represent the policy by a Gaussian distribution. i.e: $\pi_{\theta}(a|s_t) \sim N(a|f_{\theta}^{\mu}(s_t), f_{\theta}^{\Sigma}(s_t))$ [28], [29](Or a student distribution).

In PPO with KL-penalty, the surrogate objective can be seen as equation (24). The asymmetry of KL divergence will drive the policy parameter to update in the opposite direction of the π_2 parameter and will lead to the decrease of efficiency. Here we give Therom 1 to illustrate this:

Theorem 1: In the KL-PPO algorithm. When the policy's dimension is n , and $\mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] > \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right]$:

$$\min\{\beta_1, \beta_2\} \sum_{i=1}^n \log(h_i^2 e^{\frac{1-h_i^4}{2h_i^2}}) > \mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] - \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right]$$

Then $\langle \nabla_{\theta} \mathcal{L}^{KL}(\pi_2), \Delta \theta_{\pi_2|1} \rangle < 0$, Where $\frac{\sigma_{1ii}}{\sigma_{2ii}} = h_i$ represents the ratio of the i th diagonal components of the covariance matrix of π_1 , π_2 ; β_1, β_2 are the adaptive coefficients in equation (24).

Proof :

Represent the surrogate objective that update from policy $\pi_i \sim N(\mu_i, \Sigma_i)$ to $\pi_j \sim N(\mu_j, \Sigma_j)$ by $\mathcal{L}^{KL}(\pi_i|\pi_j)$. Let π_1, π_2 denote two different policy;

$$\mathcal{L}^{KL}(\pi_2|\pi_1) = \mathbb{E} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) - \beta D_{KL}(\pi_1 \parallel \pi_2) \right]$$

Similarly, we have the counterpart from π_2 to π_1 :

$$\mathcal{L}^{KL}(\pi_1|\pi_2) = \mathbb{E} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) - \beta D_{KL}(\pi_2 \parallel \pi_1) \right]$$

While $\pi_1 \neq \pi_2, D_{KL}(\pi_1 \parallel \pi_2) \neq D_{KL}(\pi_2 \parallel \pi_1)$. Then, when $\mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] > \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right]$, assume:

$$\beta_1 D_{KL}(\pi_1 \parallel \pi_2) > \beta_2 D_{KL}(\pi_2 \parallel \pi_1)$$

Therefore, we further assume this asymetry as obvious enough:

$$\beta_1 D_{KL}(\pi_1 \parallel \pi_2) - \beta_2 D_{KL}(\pi_2 \parallel \pi_1) > \mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] - \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right] \quad (27)$$

For one dimensional policy, in appendix A we have proved that:

$$\begin{aligned} & D_{KL}(\pi_1 \parallel \pi_2) - D_{KL}(\pi_2 \parallel \pi_1) \\ &= \log \frac{|\Sigma_1|}{|\Sigma_2|} + \frac{1}{2} (\text{tr}(\Sigma_1^{-1} \Sigma_2 + \Sigma_2^{-1} (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T) - \text{tr}(\Sigma_2^{-1} \Sigma_1 + \Sigma_1^{-1} (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T)) \end{aligned}$$

where $\Sigma_1^{-1} (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$ is a positive definite matrix, so the inequality is always satisfied:

$$\begin{aligned} & D_{KL}(\pi_1 \parallel \pi_2) - D_{KL}(\pi_2 \parallel \pi_1) \\ & \geq \log \frac{|\Sigma_1|}{|\Sigma_2|} + \frac{1}{2} (\text{tr}(\Sigma_1^{-1} \Sigma_2) - \text{tr}(\Sigma_2^{-1} \Sigma_1)) \end{aligned}$$

Due to the corrvarance matrix: $\Sigma = \text{diag}(\sigma_{ii}^2)$ is a diagonal matrix, we have: $\text{tr}(\Sigma_2^{-1} \Sigma_1) = \sum_{i=1}^n \frac{\sigma_{1ii}^2}{\sigma_{2ii}^2} = \sum_{i=1}^n h_i^2$ and $\frac{|\Sigma_1|}{|\Sigma_2|} = \prod_{i=1}^n h_i^2$. So we have:

$$\begin{aligned} & \beta_1 D_{KL}(\pi_1 \parallel \pi_2) - \beta_2 D_{KL}(\pi_2 \parallel \pi_1) \geq \\ & \min\{\beta_1, \beta_2\} (D_{KL}(\pi_1 \parallel \pi_2) - D_{KL}(\pi_2 \parallel \pi_1)) \geq \\ & \min\{\beta_1, \beta_2\} \sum_{i=1}^n \log(h_i^2 e^{\frac{1-h_i^4}{2h_i^2}}) \end{aligned}$$

In the inequality (27), if:

$$\min\{\beta_1, \beta_2\} \sum_{i=1}^n \log(h_i^2 e^{\frac{1-h_i^4}{2h_i^2}}) > \mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] - \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right] \quad (28)$$

Meanwhile, transpose equation 27, we have:

$$\begin{aligned} \mathcal{L}^{KL}(\pi_1|\pi_2) &= \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right] - \beta_2 D_{KL}(\pi_2 \parallel \pi_1) > \\ \mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] - \beta_1 D_{KL}(\pi_1 \parallel \pi_2) &= \mathcal{L}^{KL}(\pi_2|\pi_1) \end{aligned}$$

It illustrates the relationship of the surrogate functions:

$$\mathcal{L}^{KL}(\pi_1|\pi_2) > \mathcal{L}^{KL}(\pi_2|\pi_1) \quad (29)$$

The update of the policy is based on gradient method, we use the vector $\Delta \theta_{\pi_i|j}$ to represent the increasement of policy π_j update to policy π_i . So, for the policy update from π_1 to π_2 , we have:

$$\pi_2 \leftarrow \pi_1 + \Delta \theta_{\pi_2|1} \quad (30)$$

And for the policy gradient optimization method :

$$\Delta \theta_{\pi_2|1} \propto \nabla_{\theta} \mathcal{L}^{KL}(\pi_1) \quad (31)$$

Therefore, according to the definition of policy gradient and inequality 29, we have:

$$\langle \nabla_{\theta} \mathcal{L}^{KL}(\pi_1), \Delta \theta_{\pi_2|1} \rangle < 0 \quad (32)$$

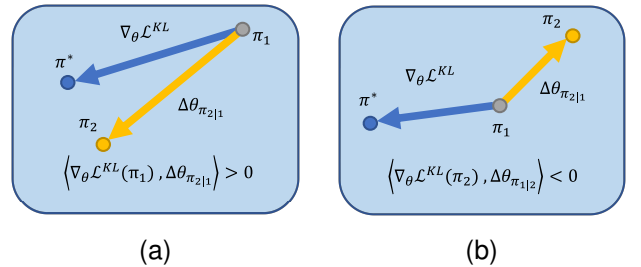


Fig. 2. The schematic diagram of policy update based on gradient. The blue arrow represent the gradient of surrogate function \mathcal{L}^{KL} on policy π_1 , the yellow arrow represent the increasement from policy π_1 to π_2 . If the intersection angle of two arrows is acute, the inner product of gradient and the increasement is positive, which indicate the new policy π^* will be draft away π_2 .

In the condition 28, although initialing we assume $\mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] > \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right]$, the surrogate functions appear to be opposite. It indicates the policy gradient method will drive policy update away π_2 , even the policy π_2 has higher $\mathbb{E}[A_{\pi}]$.

Theorem. 2 When policy's dimension is 1, let the ratio of two policy's variance $\frac{\sigma_2}{\sigma_1} = h$, while:

$$\min\{\beta_1, \beta_2\} \log(h^2 e^{\frac{1-h^4}{2h^2}}) > \mathbb{E}_{\pi_1} \left[\frac{\pi_2}{\pi_1} \hat{A}_{\pi_2}(s, a) \right] - \mathbb{E}_{\pi_2} \left[\frac{\pi_1}{\pi_2} \hat{A}_{\pi_1}(s, a) \right]$$

We can obtain the same result in Theorem 1.

The proof of Theorem 2 is obvious, we just let n in theorem equal to 1, and use a similar method on the one dimension normal distributions.

Besides, in Theorem 1, if :

$$\min\{\beta_1, \beta_2\} \log(h_i^2 e^{\frac{1-h_i^4}{2h_i^2}}) > 0, \forall i = 1, 2, \dots, n$$

then when the dimension n increases, the influence of asymmetry will increase as well. It illustrates the instability of high-dimensional space.

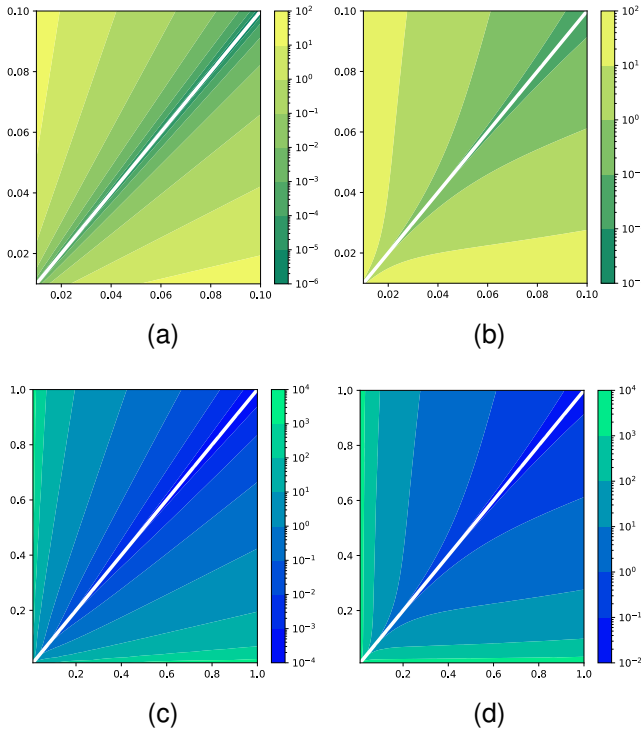


Fig. 3. Asymmetry difference of KL divergence with different parameters: (a): means: $\mu_1 = \mu_2 = 1$ and variance variant from 0.01 to 0.1; (b): means: $\mu_1 = 1, \mu_2 = 1.1$ and variance variant from 0.01 to 0.1; (c): means: $\mu_1 = 1, \mu_2 = 1.1$ and variance variant from 0.01 to 1; (d): means: $\mu_1 = 1, \mu_2 = 2$ and variance variant from 0.01 to 1 and 1 to 10. The higher value in the plots shown the stronger asymmetry in corresponding area.

In numerical, we have already given the degree of asymmetry of the KL divergence on normal distributions in appendix

A. Next, we visualize the KL divergence in Fig. 1, the figures indicate that the KL divergence changes as the feature of the policy. While changing variance in a small range, the KL divergence of two policy with the same means show slight asymmetry, as the increase of the difference of two variances, this asymmetry turns to be serious. We can see there is an order of magnitude difference between $D_{\text{KL}}(\pi_1 \parallel \pi_2)$ and $D_{\text{KL}}(\pi_2 \parallel \pi_1)$ in some area.

The KL divergence is unbounded [30], so the asymmetric KL divergence can be unlimited. Moreover, to quantitatively visualize this asymmetry, we draw the plots of the asymmetry difference of KL divergence: $\|D_{\text{KL}}(\pi_1 \parallel \pi_2) - D_{\text{KL}}(\pi_2 \parallel \pi_1)\|$ which changed by means and variance. (See Fig. 3.) From the plots, it is easy to see this asymmetry will be obvious when the means of two distributions are different, even a little different (See (b),(c),(d)). In (d), we can see the difference between $D_{\text{KL}}(\pi_1 \parallel \pi_2)$ and $D_{\text{KL}}(\pi_2 \parallel \pi_1)$ caused by asymmetry even reach four orders of magnitude:

$$\|D_{\text{KL}}(\pi_1 \parallel \pi_2) - D_{\text{KL}}(\pi_2 \parallel \pi_1)\| \sim O(10^4)$$

After theoretical analysis, we concluded several weaknesses of the KL divergence based on policy optimization algorithms:

- In TRPO, the constraint on $D_{\text{KL}}^{\text{max}}(\theta_{\text{old}}, \theta)$ is hard to numerical approximate, although it was replaced by $\overline{D}_{\text{KL}}^{\rho_{\theta_{\text{old}}}}(\theta_{\text{old}}, \theta)$ and took some numerical technique, a lot of calculations still be left.
- When TRPO meets the large-scale tasks, even the ℓ_2 norm is considered as the replacement of constraint. ℓ_2 norm is not a suitable metric to evaluate two distributions.
- In PPO with Adaptive KL Penalty, it is difficult to figure out how large steps should parameters take. In other words, the adjustment of coefficient β is not always perfect.
- Both the two algorithms introduced KL-divergence as the metric between the new policy and old policy. Due to the asymmetric of KL-divergence, it may cause vibration, especially in discrete action space.

IV. PROPOSED CIM-PPO

This section discusses the metric theorem, introduces a correntropy induced metric for probability distributions' distance evaluating, and then applies it in PPO.

A. Correntropy

ℓ_2 norm is one of the most popular measures of two vectors' similarity, due to its convenience for computation. It is always used as the approximation of the complex metric function. i.e: TRPO attempted to replace KL-divergence by ℓ_2 norm. Besides, the difference of two vectors: $(x-y)$, its inner product is equivalent to the square of its ℓ_2 norm. On the one hand, ℓ_2 norm (or two vectors difference's inner product) is a good measure because it is simple and fulfills abundant geometry means. On the other hand, it can not always represent the similarity of different data.

Kernel method is a kind of technique in statistical learning, which replaces the inner products with kernel functions that

satisfy Mercer's Theorem [31]. Any kernel function corresponds to a reproducing kernel Hilbert space (RKHS) [32] prove to see in the monograph [33]. When using a kernel function to replace the inner product, a nonlinear mapping Φ had to be done [34], which mapping the original data from the input space to an infinite-dimensional RKHS. Defined the x, y is the data in the input space. RKHS \mathbf{F} . Then, the corresponding kernel function is defined as:

$$\kappa_\sigma(x - y) = \langle \Phi(x), \Phi(y) \rangle_{\mathbf{F}} \quad (33)$$

Usually, the relationship between two data is not always simple and obvious. Using the kernel method can make the relationship more clear. e.g. Kernel Support Vector Machine(KSVM) the monograph [33], which transforms the linear inseparable data to a nonlinear space, such that the data is linearly separable.

It is interesting to define the two data's similarities in the new space. Towards this direction, Liu et.al proposed a similarity measure called **Correntropy** [35]. Since it was brought up, it was widely used to deal with non-Gaussian noise and impulsive noise in signal processing. As the development of information theoretic learning(ITL) [36], it was applied in ITL as an informational metric, to improve the performance [37]. In essence, the theory of correntropy is about noise treatment and an evaluation of the distance between two probability distributions. It is similar to Mahalanobis distance, which is a fundamental measure of two random vectors or the distance between a random vector and its center of distribution [38].

Definition: According to Renyi's quadratic entropy, correntropy is a generalized similarity measure between two arbitrary scalar random variables x and y defined as:

$$V_\sigma(x, y) = \mathbb{E} [\kappa_\sigma(x - y)] \quad (34)$$

where $\kappa_\sigma(\cdot)$ is a kernel function that satisfies Mercer's Theorem. The properties of correntropy vary with the kernel function, a good kernel function will bring a good result. In the next part, we will discuss how to choose a good kernel to gain improvement in our tasks. The common kernel functions can be seen in Tab. I. Besides, the statistical properties of correntropy can be seen in [39].

TABLE I
SEVERAL KERNEL FUNCTIONS [40]

Kernel Name	$\kappa_\sigma(x - y)$
Epanechnikov	$\min\{\frac{3}{4\sqrt{5}}(1 - \frac{1}{5}\frac{\ x-y\ ^2}{\sigma^2}), 0\}$
Biweight	$\min\{\frac{15}{16}(1 - \frac{\ x-y\ ^2}{\sigma^2})^2, 0\}$
Triangular	$\min\{1 - \frac{\ x-y\ }{\sigma}, 0\}$
Laplace	$e^{-\frac{\ x-y\ }{\sigma}}$
Gaussian	$e^{-\frac{\ x-y\ ^2}{2\sigma^2}}$
Rectangular	$\frac{1}{2}$ for $\ x - y\ < 1, 0$ otherwise

In signal processing, the most common kernel is Gaussian kernel, e.g: [41], [42]. If we use Taylor expansion onto Gaussian kernel function, we have:

$$\begin{aligned} & \exp\{-\frac{\|x-y\|^2}{2\sigma^2}\} \\ &= 1 - \frac{(x-y)^2}{2\sigma^2} + \frac{1}{2!}(\frac{(x-y)^2}{2\sigma^2})^2 - \dots + \frac{(-1)^n}{n!}(\frac{(x-y)^2}{2\sigma^2})^n \\ &= \sum_{n=0}^{\infty} \frac{(-1)^n}{n!}(\frac{(x-y)^2}{2\sigma^2})^n \end{aligned} \quad (35)$$

It reflects that Gaussian kernel is consistent of the weighted even moments, the i th weight is: $\frac{(-1)^i}{i!\sigma^{2i}}$. Easy to see, compare to high-order moments, the low-order moments have larger weight, and will occupy a more significant position. As the increase of σ , the value of higher-order moments decreases slightly. In a view of information, the size of σ will affect the information content of high-order moments in correntropy. Conventionally, the kernel sizes or bandwidth is a free parameter that must be chosen by the user using concepts of density estimation, such as Silverman's rule [40] or maximum likelihood. Although the Gaussian kernel shows its great character, in the view of computing, the Gaussian kernel is equal to compute n times all the quadratic kernel, in RL tasks, it is not a computation-friendly kernel. Therefore, in the latter part, we will utilize other kernel functions to improve.

B. CIM-PPO

In the previous part, we discussed the weakness of the present method in policy optimization and the concept of correntropy. This section will discuss how to improve the algorithm by introducing correntropy.

In statistics view, KL divergence is a distance function between two probability distributions and is considered as a kind of metric. However, KL divergence is not a metric, because it does not profile the properties [43].

Definition: A function $d: X \times X \rightarrow \mathbb{R} \cup \{\infty\}$ is a metric on the set X if for all $x, y, z \in X$, the three follow properties are satisfying [43]:

- 1) **Positiveness:** $d(x, y) \geq 0$, $d(x, y) = 0$ iff $x = y$.
- 2) **Symmetry:** $d(x, y) = d(y, x)$
- 3) **Triangle inequality:** $d(x, z) \leq d(x, y) + d(y, z)$.

The concept of "Metric" is similar to "distance" in geometry. It is widely used as the measurement of two things' differences. For example, the ℓ_2 norm is a well-known metric on euclidean space. Therefore, KL-divergence is not a metric, because it is an asymmetric function: $D_{KL}(x, y) \neq D_{KL}(y, x)$. Correntropy is not a metric as well, due to $V_\sigma(x, y) \neq 0$, while $x = y$. For extending correntropy to metric, Liu et al [35] further proposed **Correntropy Induced Metric (CIM)**, which can be expressed as follow:

$$CIM(x, y) = (V_\sigma(0) - V_\sigma(x, y))^{\frac{1}{2}} \quad (36)$$

CIM has several good properties:

Property 1: CIM is positive and can be bounded.

It is easy to see the minimum value of CIM is 0 while $x - y = 0$ by equation (36), and then, for example, if we chose Gaussian function(Tab. I) as the kernel function of correntropy, using the singularity of Gaussian kernel, we can

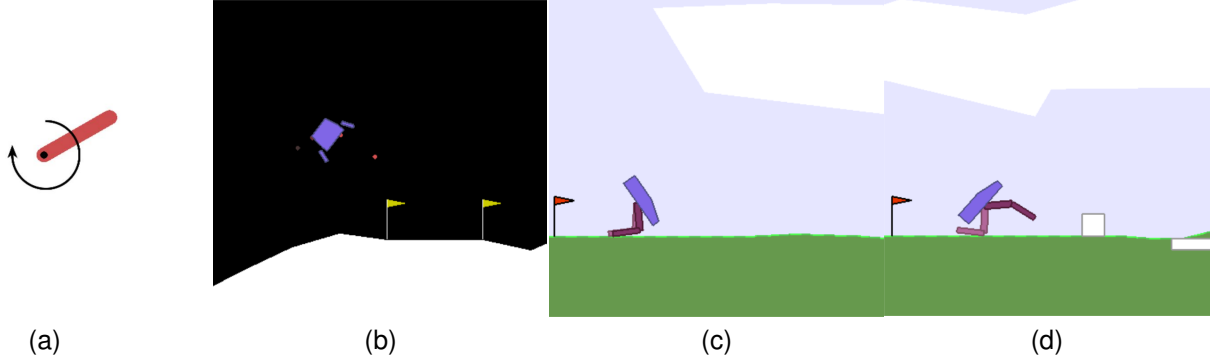


Fig. 4. Screenshots of the environments for all testing tasks. All of these are continuous tasks. (a) Pendulum. (b) LunarLanderContinuous. (c) BipedalWalker. (d) BipedalWalkerHardcore.

obtain the superior of CIM, then the inequality is satisfied: $0 \leq CIM_\sigma(x, y) < 1$. This property indicates the CIM can be bounded while the KL divergence is unbounded.

Property 2: CIM can be a symmetric function: $CIM_\sigma(x, y) = CIM_\sigma(y, x)$

When the kernel function of correntropy is symmetric $\mathbb{E}[\kappa_\sigma(x - y)] = \mathbb{E}[\kappa_\sigma(y - x)]$, the correntropy will be symmetric: $V_\sigma(x, y) = V_\sigma(y, x)$

Therefore, we have:

$$\begin{aligned} CIM(x, y) &= (V_\sigma(0) - V_\sigma(x, y))^{\frac{1}{2}} \\ &= (V_\sigma(0) - V_\sigma(y, x))^{\frac{1}{2}} = CIM(y, x) \end{aligned}$$

Property 3: For any kernel function, CIM satisfies the triangle inequality: $CIM_\sigma(x, y) \leq CIM_\sigma(x, z) + CIM_\sigma(z, y)$. The proof see in [35] or appendix B.

The above 3 properties indicate that CIM can be a metric while the kernel function is symmetric, and compared to the unbounded KL divergence, CIM is always bounded, and more stationary than KL divergence. Moreover, CIM has pretty good robustness, especially in two distributions distance evaluation. It can deal with not only the distributions with Gaussian noise but also with non-gaussian noise [35]. In TRPO and KL-PPO, KL-divergence is used as the measure of the difference between the new policy and the old policy. In essence, the policy is a distribution, and the KL-divergence is a measure of different distributions, as well as CIM. However, the asymmetry and large-scale calculation make the policy optimization method performance not always perfect. Therefore, we use the correntropy mechanism to replace the adaptive parameter adjustment mechanism in KL-PPO, and the CIM will substitute the KL divergence as the metric to evaluate the distance between the new policy and the old policy. And then, for eliminating the effect of an asymmetric penalty and reducing its calculation in KL-PPO, the computationally expensive kernel functions will not be selected.

To sum up, we modified Proximal Policy Optimization with KL divergence and proposed PPO with objective Correntropy Induced Metric (CIM-PPO), the surrogate of CIM-PPO is present here:

$$\begin{aligned} \mathcal{L}_\pi^{CIM} &= \\ \mathbb{E}_\pi &\left[\frac{\pi_{new}(a|s)}{\pi(a|s)} \hat{A}_\pi(s, a) - \alpha CIM_\sigma(\pi(\cdot|s), \pi_{new}(\cdot|s)) \right] \end{aligned}$$

Algorithm 2 PPO with Correntropy Induced Metric

Input: Initial policy: π_0 . Initial Actor-Critic parameter ϕ_0 .

1. According to the task, select penalty control parameter α
2. Choice whether estimate σ or set as 1.

for $i = 0, 1, 2, \dots$ **until convergence do**

3. Running policy π_i in the task environment, and storage the trajectories in the set $\mathcal{D}_i = \{\tau_k\}$.
4. Estimate the reward \hat{R}_t .
5. Estimate \hat{A}_t by Actor-Critic based on current parameter ϕ_k .
6. Set $\sigma = 1$ or estimate $\hat{\sigma}$ by Mercer's law
7. Optimize follow estimation by minibatch SGD or Adam:

$$\pi_{i+1} = \arg \max_{\pi} \frac{1}{|D_i| \cdot N} \sum_{\tau \in D_i} \sum_{k=0}^N \frac{\pi}{\pi_i} \hat{A}_t - \alpha CIM_{\hat{\sigma}}(\pi_i, \pi)$$

8. Update ϕ_k through value function V_ϕ of Actor-Critic:

$$\phi_{i+1} = \arg \min_{\phi} \frac{1}{|D_i| \cdot N} \sum_{\tau \in D_i} \sum_{k=0}^N (V_\phi(s_t) - \hat{R}_t)^2$$

end for

Where α is the constant that is based on the task, if we do not hope the difference between the new policy and old policy is too large in each update, we can set α a large number. If we don't consider it or we permit two policies to have differences in a range, we can set α as a smaller number. The surrogate objective is optimized by stochastic gradient descent (SGD) or Adam [44]. Algorithm 2 is the pseudocode of CIM-PPO.

Compared with KL-PPO, CIM-PPO uses correntropy induced metric as the penalty in the surrogate objective. Meanwhile, the adaptive adjustment mechanism was canceled in CIM-PPO due to enough robustness of CIM. To validate the performance of CIM-PPO, in the next section, we will give the experiments.

Overall, we design experiments to compare the efficiency of our algorithm with Clip-PPO and KL-PPO. Choose four basic

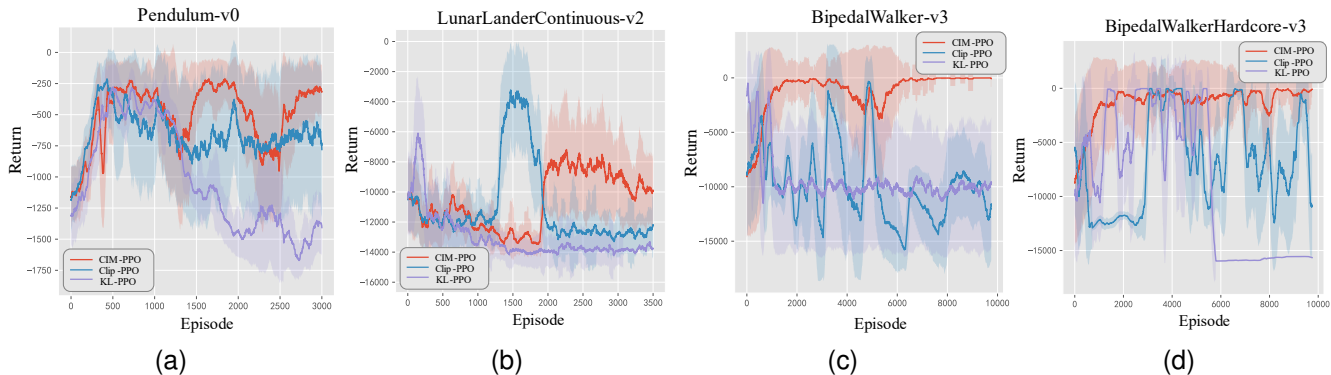


Fig. 5. Experiment results: CIM-PPO and two classic PPO algorithm’s learning curves on four basic Gym environments: (a)Pendulum-v0. (b)LunarLanderContinuous-v2. (c) BipedalWalker-v3. (d) BipedalWalkerHardcore-v3.

TABLE II
THE INFORMATION OF THE FOUR TASKS

Task name	Action Space dimension	State Space dimension
Pendulum-v0	1	3
LunarLanderContinuous-v2	2	8
BipedalWalker-v3	4	24
BipedalWalkerHardcore-v3	4	24

TABLE III
THE HYPERPARAMETERS SETTING

Algorithms name	Hyperparameters name	Value
KL-PPO	d_{target}	0.1
KL-PPO	Initial β	0.5
Clip-PPO	ϵ	0.2
CIM-PPO	α	1
KL/Clip/CIM-PPO	γ	0.9
KL/Clip/CIM-PPO	Learning rate of Actor networks	0.0001
KL/Clip/CIM-PPO	Learning rate of Critic networks	0.0002
KL/Clip/CIM-PPO	Update batchsize	32
KL/Clip/CIM-PPO	Critic update steps	10
KL/Clip/CIM-PPO	Actor update steps	10

continuous tasks from OpenAI gym(See Fig. 4 and Tab II). In the four tasks, the dimension of action space increases from 1 to 4, and the corresponding state space increases from 3 to 24. The normal distributions’ variances of these four tasks in the construction of policy vary from 0.1 to 1, which will make the asymmetry of KL divergence more obvious. In theorem 1, we illustrated that as the dimension increases, the asymmetry of KL penalty will be obvious, the performance of KL-PPO will decline.

In the selection of kernel functions, based on the principle of Occam’s razor and the previous discussion in the part of correntropy, we choose a simple kernel. So the Triangular kernel was chosen as the first two tasks’ kernel function, the Biweight kernel was chosen as the rest tasks, then, training agents in four basic tasks. The hyperparameters setting of CIM-PPO, KL-PPO, and Clip-PPO can be seen in Tab. III

V. RESULTS

To evaluate the performance of CIM-PPO, we compute the learning curves of four tasks, and the results are shown in Fig. 5.

Overall, in the learning rate, CIM-PPO can at least reach the equal effect to Clip-PPO and sometimes better than Clip-PPO, always perform better than KL-PPO.

Specifically, in the Pendulum task, initially, all three algorithms can achieve the almost equal effect, but KL-PPO collapsed in about 1200 episodes, Clip-PPO tends to be stationary, although CIM-PPO vibrates in a range, the reward is still higher than Clip-PPO. In LunarLanderContinuous task, the Clip-PPO and KL-PPO all experience the process that rewards spike at first but quickly decreases to the original value, but CIM-PPO increase and keep its reward at a high level for a long time, which indicate that CIM-PPO has better robustness.

Our algorithm in the last two tasks performed best. In Fig. 5(c),(d), we can see the reward of CIM-PPO steadily rise with the training episodes, before it reaches its limitation.

VI. CONCLUSION

PPO is one of the most famous algorithms in deep reinforcement learning, it performs excellently in many challenging tasks. However, the PPO is not always perfect, for example,

Clip-PPO lacks theoretical explanation in clip operation; KL-PPO is not robust.

Our research mainly focuses on KL-PPO, we studied the asymmetry of KL divergence in KL-PPO and proved this asymmetry may affect the robustness and learning efficiency. We proved this influence of KL-divergence will be obvious as the increase of the policy's dimension. We have introduced the Correntropy mechanism into KL-PPO to replace the penalty mechanism, and use CIM (a robust metric that was widely used in M-estimation to evaluate two distributions' difference) to substitute the KL divergence as the measure of distance between the old policy and new policy. In the view of the signal process, by sampling the task's environment and using parameterized tricks, the policy is constructed as a normal distribution or a student distribution. It is hard to ensure the sampling does not exist non-gaussian noise, but the CIM performance is well and robust in dealing with non-gaussian noise.

Shortly, we may extend the proposed results along with two interesting directions. The one is to deeply analyze the reason why CIM can make PPO performs better, and put forward a complete theory to guide researchers choose the best kernel function in a specific task. The other one is to study the effect of kernel size σ and give an optimal kernel size estimate method that can improve the tasks' performance most.

REFERENCES

- [1] X. Wang, T. Li, and Y. Cheng, "Proximal parameter distribution optimization," *IEEE Trans. Syst., Man, Cybern., Syst.*, early access, Aug. 30, 2019, DOI: 10.1109/TSMC.2019.2931946.
- [2] Y. Gu, Y. Cheng, C. L. P. Chen and X. Wang, "Proximal Policy Optimization With Policy Feedback," in *IEEE Trans. Syst., Man, Cybern., Syst.*, DOI: 10.1109/TSMC.2021.3098451.
- [3] X. Wang, T. Li, Y. Cheng, and C. L. P. Chen, "Inference-based posterior parameter distribution optimization," *IEEE Trans. Cybern.*, early access, Oct. 7, 2020, DOI: 10.1109/TCYB.2020.3023127.
- [4] P. Lv, X. Wang, Y. Cheng, Z. Duan, and C. L. P. Chen, "Integrated double estimator architecture for reinforcement learning," *IEEE Trans. Cybern.*, early access, Oct. 7, 2020, DOI: 10.1109/TCYB.2020.3023033.
- [5] Mnih, V., Kavukcuoglu, K., Silver, D. et al. Human-level control through deep reinforcement learning. *Nature* 518, 529–533 (2015).
- [6] Silver D, Huang A, Maddison CJ, et al. Mastering the game of Go with deep neural networks and tree search. *Nature*. 2016;529(7587):484-489.
- [7] Hasselt H V, Guez A, Silver D . Deep Reinforcement Learning with Double Q-learning. *Computer science*, 2015.
- [8] Wang Z, Schaul T, Hessel M, van Hasselt H, Lanctot M, de Freitas N. Dueling Network Architectures for Deep Reinforcement Learning.
- [9] H. Hu, S. Song, and G. Huang, "Self-attention-based temporary curiosity in reinforcement learning exploration," *IEEE Trans. Syst., Man, Cybern., Syst.*, early access, Dec. 18, 2019, DOI: 10.1109/TSMC.2019.2957051
- [10] J. Yu, Z. Wu, X. Yang, Y. Yang, and P. Zhang, "Underwater target tracking control of an untethered robotic fish with a camera stabilizer," *IEEE Trans. Syst., Man, Cybern., Syst.*, early access, Jan. 13, 2020, DOI: 10.1109/TSMC.2019.2963246.
- [11] X. You, X. Li, Y. Xu, H. Feng, and J. Zhao, "Toward packet routing with fully distributed multiagent deep reinforcement learning," *IEEE Trans. Syst., Man, Cybern., Syst.*, early access, Aug. 13, 2020, DOI: 10.1109/TSMC.2020.3012832.
- [12] X. Xu, L. Zuo, X. Li, L. L. Qian, J. K. Ren, and Z. P. Sun, "A reinforcement learning approach to autonomous decision making of intelligent vehicles on highways," *IEEE Trans. Syst., Man, Cybern., Syst.*, vol. 50, no. 10, pp. 3884–3897, Oct. 2020.
- [13] Y. Keneshloo, T. Shi, N. Ramakrishnan, and C. K. Reddy, "Deep reinforcement learning for sequence-to-sequence models," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 31, no. 7, pp. 2469–2489, Jul. 2020.
- [14] Volodymyr Mnih et. al. Asynchronous Methods for Deep Reinforcement Learning, 2016, arXiv 1602.01783 [cs.LG]
- [15] Tuomas Haarnoja et. al. Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor, arXiv 1801.01290[cs.LG]
- [16] Kakade S M . A Natural Policy Gradient. *advances in neural information processing systems*, 2001.
- [17] Kakade, Sham and Langford, John. Approximately optimal approximate reinforcement learning. In *ICML*, volume 2, pp. 267–274, 2002.
- [18] Levine, Sergey and Abbeel, Pieter. Learning neural network policies with guided policy search under unknown dynamics. In *Advances in Neural Information Processing Systems*, pp.1071–1079, 2014.
- [19] Achiam J , Held D , Tamar A , et al. Constrained Policy Optimization. 2017, arXiv1705.10528 [cs.LG]
- [20] John Schulman, Filip Wolski, and Prafulla Dhariwal et. al. Proximal Policy Optimization Algorithms. 2017, arXiv 1707.06347 [cs.LG]
- [21] Taisuke Kobayashi, "Proximal Policy Optimization with Relative Pearson Divergence", arXiv : 2010.03290v2 [cs.LG] .
- [22] Y. Cheng, L. Huang and X. Wang, "Authentic Boundary Proximal Policy Optimization," in *IEEE Transactions on Cybernetics*, DOI: 10.1109
- [23] R.Sutton, A.Barto. Reinforcement Learning: An Introduction. The MIT Press Massachusetts. 2017. 62-75.
- [24] M. Puterman. Markov decision processes: Discrete stochastic dynamic programming. John Wiley and Sons, 1994.
- [25] A Agarwal, N Jiang, S Kakade et al. Reinforcement Learning: Theory and Algorithms. 2020.
- [26] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel. Trust region policy optimization. In: *CoRR*, abs/1502.05477 (2015).
- [27] Y. Wang, H. He, and X. Tan. "Truly proximal policy optimization," in *Uncertainty in Artificial Intelligence*. PMLR, 2020, pp. 113–122.
- [28] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- [29] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning (ICML)*, pages 1928–1937, 2016.
- [30] Nielsen F. On the Jensen–Shannon Symmetrization of Distances Relying on Abstract Means. *Entropy*. 2019;21(5):485.
- [31] V. Vapnik, *The Nature of Statistical Learning Theory*. New York: Springer-Verlag, 1995.
- [32] Jian-Wu Xu, A. R. C. Paiva, I. Park, and J. C. Principe, "A Reproducing Kernel Hilbert Space Framework for Information-Theoretic Learning," in *IEEE Transactions on Signal Processing*, vol. 56, no. 12, pp. 5891–5902, Dec. 2008
- [33] Schölkopf, Bernhard. Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT Press, 2003.
- [34] Wang L, Solomon J, Gehre A, "Bronstein MM. Kernel Functional Maps. *Computer Graphics Forum*." 2018;37(5):27-36.
- [35] W. Liu, P. P. Pokharel, and J. C. Principe, "Correntropy: Properties and applications in non-Gaussian signal processing," *IEEE Trans. Signal Process.*, vol. 55, no. 11, pp. 5286–5298, Nov. 2007
- [36] D. Erdogmus and J. C. Principe, "An error-entropy minimization algorithm for supervised training of nonlinear adaptive systems," *IEEE Trans. Signal Process.* vol. 50, pp. 1780–1786, Jul. 2002.
- [37] R. He, Y. Zhang, Z. Sun and Q. Yin, "Robust Subspace Clustering With Complex Noise," in *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 4001-4013, Nov. 2015
- [38] Dai D, Liang Y. High-Dimensional Mahalanobis Distances of Complex Random Vectors. *Mathematics*. 2021; 9(16):1877.
- [39] I.Santamaria, P.P.Pokharel, and J.C.Principe, "Generalized correlation function: Definition, properties, and application to blind equalization," *IEEE Trans. Signal Process.*, vol.54, no.6, pp.2187-2197, Jun. 2006.
- [40] B. W. Silverman, *Density Estimation for Statistics and Data Analysis*. London, U.K.: Chapman and Hall, 1986.
- [41] Guobing Qian, Fuliang He, Shiyuan Wang, H.C. Iu Herbert, Robust constrained maximum total correntropy algorithm, *Signal Processing*, Volume 181, 2021, 107903,
- [42] L. Liu, Q. Feng, C.L.P. Chen and Y. Wang, "Noise Robust Face Hallucination Based on Smooth Correntropy Representation," in *IEEE Transactions on Neural Networks and Learning Systems*, DOI: 10.1109/TNNLS.2021.3071982.
- [43] D. Burago, Y Burago. & S. Ivanov. A Course in Metric Geometry. American Mathematical Society.
- [44] D. Kingma and J. Ba. "Adam: A method for stochastic optimization". arXiv:1412.6980 (2014).

APPENDIX

A. The asymmetry of KL divergence

1) **Case 1:** In one dimensional normal distribution:

Let $p \sim N(\mu_1, \sigma_1)$, $q \sim N(\mu_2, \sigma_2)$, The KL divergence of p and q can be given by:

$$\begin{aligned} D_{\text{KL}}(p \parallel q) &= \int_{-\infty}^{+\infty} p(x) \log \frac{p(x)}{q(x)} dx \\ &= - \int_{-\infty}^{+\infty} p(x) \log[q(x)] dx + \int_{-\infty}^{+\infty} p(x) \log[p(x)] dx \end{aligned}$$

To derive the expression of KL divergence, we consider the first part of formula above:

$$\begin{aligned} &- \int_{-\infty}^{+\infty} p(x) \log[q(x)] dx \\ &= - \int_{-\infty}^{+\infty} p(x) \log \left[\frac{1}{\sqrt{2\pi\sigma_2^2}} \exp\left\{-\frac{(x-\mu_2)^2}{2\sigma_2^2}\right\} \right] dx \\ &= \int_{-\infty}^{+\infty} p(x) \log \sqrt{2\pi\sigma_2^2} dx + \int_{-\infty}^{+\infty} p(x) \frac{(x-\mu_2)^2}{2\sigma_2^2} dx \\ &= \int_{-\infty}^{+\infty} p(x) \log \sqrt{2\pi\sigma_2^2} dx + \\ &\quad \frac{1}{2\sigma_2^2} \int_{-\infty}^{+\infty} p(x) (x^2 - 2\mu_2 x + \mu_2^2) dx \end{aligned}$$

By the properties of a random variable:

$$\begin{aligned} \int_{-\infty}^{+\infty} a \cdot p(x) dx &= a, \int_{-\infty}^{+\infty} ax \cdot p(x) dx = a \cdot \mathbb{E}_p[x] \\ \int_{-\infty}^{+\infty} ax^2 \cdot p(x) dx &= a \cdot ((\mathbb{E}_p[x])^2 + \text{Var}[x]) \end{aligned}$$

We have:

$$\begin{aligned} &\int_{-\infty}^{+\infty} p(x) \log \sqrt{2\pi\sigma_2^2} dx + \frac{1}{2\sigma_2^2} \int_{-\infty}^{+\infty} p(x) (x^2 - 2\mu_2 x + \\ &\mu_2^2) dx \\ &= \sqrt{2\pi\sigma_2^2} + \frac{\mu_1^2 + \sigma_1^2}{2\sigma_2^2} - \frac{2\mu_1\mu_2}{2\sigma_2^2} + \frac{\mu_2^2}{2\sigma_2^2} \\ &= \log \sqrt{2\pi\sigma_2} + \frac{(\mu_2 - \mu_1)^2 + \sigma_1^2}{2\sigma_2^2} \end{aligned}$$

By using the expression of the first part, the last part of KL divergence formula can be derived:

$$\begin{aligned} &\int_{-\infty}^{+\infty} p(x) \log[p(x)] dx \\ &= -\log \sqrt{2\pi\sigma_1} - \frac{\sigma_1^2}{2\sigma_1^2} \end{aligned}$$

Combining the two parts expression, we can obtain the expression of KL divergence with one dimensional normal distribution:

$$D_{\text{KL}}(p \parallel q) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

Similarly, we have:

$$D_{\text{KL}}(q \parallel p) = \log \frac{\sigma_1}{\sigma_2} + \frac{\sigma_2^2 + (\mu_2 - \mu_1)^2}{2\sigma_1^2} - \frac{1}{2}$$

When $\sigma_1 \neq \sigma_2$, $D_{\text{KL}}(p \parallel q) - D_{\text{KL}}(q \parallel p) \neq 0$.

It is easy to see the KL divergence is asymmetry when the distributions are normal. Moreover, we define the asymmetry difference of KL divergence for evaluating the asymmetry numerically:

$$\begin{aligned} D_{\text{KL}}(p \parallel q) - D_{\text{KL}}(q \parallel p) &= \log \left(\frac{\sigma_2}{\sigma_1} \right)^2 \\ &+ \frac{(\sigma_1^2 - \sigma_2^2)[(\mu_1 - \mu_2)^2 + (\sigma_1^2 + \sigma_2^2)]}{2\sigma_1^2\sigma_2^2} \end{aligned}$$

2) **Case 2:** In multi-dimensional normal distribution:

Consider two multi-dimensional normal distribution: $\mathbf{P} \sim N(\mu_1, \Sigma_1)$, $\mathbf{Q} \sim N(\mu_2, \Sigma_2)$. The KL divergence between P and Q is:

$$\begin{aligned} D_{\text{KL}}(P \parallel Q) &= \int P(x) \log \frac{P(x)}{Q(x)} dx \\ &= - \int P(x) \log[Q(x)] dx + \int P(x) \log[P(x)] dx \end{aligned}$$

Considering $\int P(x) \log[Q(x)] dx$, expanding $Q(x)$ by multi-normal distribution's expression:

$$\begin{aligned} &\int P(x) \log[Q(x)] dx \\ &= \int P(x) \log[(2\pi)^{-\frac{n}{2}} |\Sigma_2|^{\frac{1}{2}} \exp\left\{-\frac{1}{2}(x-\mu_2)^T \Sigma_2^{-1} \right. \\ &(x-\mu_2)\left.\right\}] dx \\ &= \int P(x) \left(-\frac{n}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2|\right) dx - \\ &\frac{1}{2} \int (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2) P(x) dx \\ &= -\frac{1}{2} \int (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2) P(x) dx \\ &- \frac{n}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2| \end{aligned}$$

Because the dimension of $(x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)$ is one, so we have $\text{tr}((x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)) = (x-\mu_2)^T \Sigma_2^{-1} (x-\mu_2)$. Due to: $\text{tr}(ABC) = \text{tr}(CBA)$, if ABC and CBA existed. So we have:

$$\begin{aligned}
& \frac{1}{2} \int (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) P(x) dx \\
&= \frac{1}{2} \int \text{tr}((x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2) P(x)) dx \\
&= \frac{1}{2} \int \text{tr}((x - \mu_2)(x - \mu_2)^T \Sigma_2^{-1} P(x)) dx \\
&= \frac{1}{2} \int \text{tr}([(x - \mu_1) + (\mu_1 - \mu_2)][(x - \mu_1) + (\mu_1 - \mu_2)]^T \Sigma_2^{-1} P(x)) dx \\
&= \frac{1}{2} \left(\int \text{tr}((x - \mu_1)(x - \mu_1)^T \Sigma_2^{-1} P(x)) dx \right. \\
&\quad \left. - 2 \int \text{tr}((\mu_2 - \mu_1)(x - \mu_1)^T \Sigma_2^{-1} P(x)) dx + \int \text{tr}(\Sigma_2^{-1} (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T P(x)) dx \right) \\
&= -\frac{1}{2} \text{tr}(\Sigma_2^{-1} \Sigma_1) - \frac{1}{2} \text{tr}(\Sigma_2^{-1} (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T)
\end{aligned}$$

Therefore, we obtain the expression as follows:

$$\begin{aligned}
& \int P(x) \log[Q(x)] dx = -\frac{n}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2| + \\
& \quad \frac{1}{2} \text{tr}(\Sigma_2^{-1} \Sigma_1) + \frac{1}{2} \text{tr}(\Sigma_2^{-1} (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T)
\end{aligned}$$

By the same way, we have:

$$\begin{aligned}
& \int P(x) \log[P(x)] dx \\
&= -\frac{n}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_1| - \frac{n}{2}
\end{aligned}$$

Finally, we have the KL divergence of multi-dimensional normal distribution :

$$\begin{aligned}
D_{\text{KL}}(P \parallel Q) &= \int P(x) \log \frac{P(x)}{Q(x)} dx \\
&= - \int P(x) \log[Q(x)] dx + \int P(x) \log[P(x)] dx \\
&= \frac{1}{2} \log \frac{|\Sigma_1|}{|\Sigma_2|} + \frac{1}{2} \text{tr}(\Sigma_1^{-1} \Sigma_2 + \Sigma_2^{-1} (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T) - \frac{n}{2}
\end{aligned}$$

It is similar to the result of KL divergence in one dimension, the $D_{\text{KL}}(P \parallel Q) \neq D_{\text{KL}}(Q \parallel P)$:

$$\begin{aligned}
& D_{\text{KL}}(P \parallel Q) - D_{\text{KL}}(Q \parallel P) \\
&= \log \frac{|\Sigma_1|}{|\Sigma_2|} + \frac{1}{2} (\text{tr}(\Sigma_1^{-1} \Sigma_2 + \Sigma_2^{-1} (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T) - \\
& \quad \text{tr}(\Sigma_2^{-1} \Sigma_1 + \Sigma_1^{-1} (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T))
\end{aligned}$$

B. The proof of CIM's property 3

Denote $X, Y, Z \in \mathbb{R}^n$. In a well defined Hilbert space \mathbf{F}^n , construct two features: $X', Y' \in \mathbf{F}^n$, which are mapped X, Y separately. By the definition of kernel function (equation (33)), the ℓ_2 norm in $X' - Y'$ can be represented as follow:

$$\begin{aligned}
\|X' - Y'\|_2 &= (\langle X' - Y', X' - Y' \rangle)^{\frac{1}{2}} \\
&= (\langle X', X' \rangle + \langle Y', Y' \rangle - 2 \langle X', Y' \rangle)^{\frac{1}{2}} \\
&= (\kappa(X - X) + \kappa(Y - Y) - 2\kappa(X - Y))^{\frac{1}{2}} \\
&= (2\kappa(0) - 2\kappa(X - Y))^{\frac{1}{2}} \\
&= \text{CIM}(X, Y)
\end{aligned}$$

ℓ_2 norm satisfies triangle inequality: $\|X - Z\|_2 \leq \|X - Y\|_2 + \|Y - Z\|_2$, we have:

$$\begin{aligned}
\text{CIM}(X, Z) &= \|X' - Z'\|_2 \\
\|X' - Y'\|_2 + \|Y' - Z'\|_2 &= \text{CIM}(X, Y) + \text{CIM}(Y, Z) \\
\text{CIM}(X, Z) &\leq \text{CIM}(X, Y) + \text{CIM}(Y, Z)
\end{aligned}$$