
Persistence Length-Based Exploration in Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 BLABLA

2 1 Introduction

3 We propose exploration algorithm that not only consider the temperature, it also considers the general
4 form of the result trajectory induced from a certain l_p .

5 We also emphasise on the fact that local self-avoidance is achieved only w.r.t. the action-space.

6 We also emphasize on the comparison between Monte-Carlo sampling and tree backup model with
7 the proposed trajectory based exploration. Since, using L_p we explore trajectories as opposed to the
8 state-action pairs individually.

9 We emphasise on trajectory-based exploration rather than the state-based exploration as well as a
10 technique that relies on local information about visited state-action pairs as opposed to relying on
11 the statistical summery of visited state-action pairs. The exploration method employs a range of
12 persistence lengths and and and explore a range of low-dimensional manifolds and then explore the
13 selected manifold efficiently. Intuitively we need to show that by choosing proper values for d_β one
14 can explore lower dimensional manifolds efficiently and then use results in IJCAI paper and explore
15 the 3D manifold. Therefore two main questions arise here:

16 1- How to choose d_B to land of the low dimensional manifold, given that such manifold exists.

17 2- If we assume that we have landed on the 3D manifold are we still going to be able to use results
18 from IJCAI and show that they still hold.

19 2 The notion of persistence length: formal treatment

20 *Maziar:* Mathematically define what a *stiff chain* is.

21 *Maziar:* Define what end-to-end vector is a stiff chain and show that the end-to-end distribu-
tion of a stiff chain is always different than the *freely-jointed* one.

22 *Maziar:* Computationally show that how the stiffness diminishes over time and change of
dimensionality

23 Note that the strong interaction between bonds reduces the effect of thermal energy.

2.1 Stiff chain, a discrete model:

A stiff chain composed of N unit vectors $\{u_1, \dots, u_N\}$, in the discrete model is parametrized by the following bending energy:

$$\mathcal{E}_B^N = \frac{\kappa}{2a} \sum_{i=1}^N (u_{i+1} - u_i)^2 \quad (1)$$

We also have the distribution function ν over the first and last bond vectors, u_f and u_l for a stiff chain of N bond vectors. Eq. 15.114.

$$\mathbb{E}_\nu[u_i \cdot u_j] = e^{\frac{-|i-j|a}{L_p}} \quad (2)$$

$$L_p = \frac{-a}{\log \left[\frac{I_{d/2}(\frac{\kappa}{a k_B T})}{I_{d/2-1}(\frac{\kappa}{a k_B T})} \right]} \quad (3)$$

2.2 Stiff chain, a continuous model:

$$\begin{aligned} \mathcal{E}_B &= \frac{\kappa}{2} \int_0^L ds (\partial_s u)^2 \\ u(s) &= \frac{d}{ds} x(s) \\ d(s) &= \sqrt{dx^2} \\ L_p &= 2 \frac{\kappa}{k_B T} \end{aligned}$$

Our intention is to guarantee the local self-avoidance within a trajectory induced by the exploration policy. We do this, firstly, by assuming that we are in the crossover regime where $L_p \sim D$. This means, that agent always feels that is encompassed within a ball of diameter $D = O(L_p)$. In the crossover regime, the statistical mechanics is governed by the competition of the *Thermal Energy* (\mathcal{E}_T) and bending rigidity κ .

Note: We have to note that even when a technique employs both action and state in the exploration phase, we still need to know the reward structure or the observation in order to decide about the next choice of actions.

Maziar: Define three different regime and see if the existing categorization extends to higher dimensional regime, clearly the intuition behind all existing definitions come from effort in identifying a parameter that describes three different regimes: 1) Weak confinement, 2) Strong confinement, 3) Crossover.

The physical quantity that captures the behavior of worm-like chain under length constraint (or confinement) is *tangent correlation function*.

In the model proposed thus far, they assume that polymer is a one-dimensional object located in a d -dimensional environment. This clearly could be problematic since in our domain trajectory/chain is not necessarily one-dimensional. Thereby, κ^I that is assumed to be invariant w.r.t. change in domain dimensionality will not invariant any more.

Maziar: We show that a chain is locally homeomorphic to Euclidean 1-space therefore it is a 1-dimensional object. The main idea comes from the fact that trajectory is a time dependent entity and time is the only free parameter and therefore we are not going to have notions like figure eight objects.

Maziar: I need to define the *internal* and *external* dimension.

Maziar: Prove where L_p^d formula is coming from.

48 Fro a chain with one internal dimension and d external dimension, L_p^d represents the persistence
 49 length in d -dimensional space. We define the intrinsic persistence length L_p^i of a chain as a property
 50 that captures the inherent competition between *Thermal Energy* and *Rigidity*.

$$L_p^i \propto \frac{\kappa^i}{\mathcal{E}_T} \quad (4)$$

51 Where κ^i represents the intrinsic rigidity required to de-correlate two adjacent bond vectors (vector
 52 that connect adjacent actions) in a chain lying on a 3-dimensional manifold and \mathcal{E}_T represents the
 53 thermal energy. One can show $L_p^d = \frac{2L_p^i}{d-d_B}$, where $d = d_T + d_B$ (or $d \geq d_T + d_B$) and d_T is the
 54 dimension of space governed by entropic energy and d_B represent the dimension of space governed
 55 by the bending energy.

56 3 Algorithm:Old

Algorithm 1: PolyRL

Input: \mathcal{A}, \mathcal{S} , Intrinsic Rigidity κ^i , temperature \mathcal{E}_T , l_B, l , Step-size $b_o, \phi, \gamma, \rho, \epsilon, \alpha, \sigma, T$

Output: \hat{w}

```

1 Construct the intrinsic  $L_p^i \leftarrow \frac{\kappa^i}{\mathcal{E}_T}$ ;
2  $L_p^l \leftarrow \frac{2L_p^i}{(l-l_B)}$ ;
3  $\mu \leftarrow \cos^{-1}(e^{-\frac{b_o}{L_p^i}})$ ;
4 Sample  $A_0$  and  $S_0$  w.r.t  $\rho$ ;
5 Let  $w_0 = 0$ ;
6 for  $t = 1$  to  $t = T$  do
7   Convert  $A_t$  to its spherical equivalent  $A_t^{sp}$ ;
8   for  $i = 1$  to  $l_\beta$  do
9     Sample  $\theta_o$  w.r.t.  $\mathcal{N}(\mu, \sigma^2)$ ;
10    Set  $A_t^{sp}(i)$  w.r.t.  $\theta_o$  and  $\epsilon$  ( $\epsilon$  greedy) //  $A_t(i)$  is the  $i$ th coordinate of  $A_t$ 
11    ;
12   for  $i = l_\beta + 1$  to  $l$  do
13     Sample  $\theta_i$  w.r.t.  $\mathcal{U}([0, 2\pi])$ ;
14     Choose  $A_t^{sp}(i)$  w.r.t.  $\theta_i$  and  $\epsilon$  ( $\epsilon$  greedy);
15   OR;
16   Convert  $A_t^{sp}$  to its Cartesian equivalent  $A_t$ ;
17   Take action  $A_t$  and observe  $R_{t+1}$  and  $S_{t+1}$ ;
18   Update  $\hat{w}$  with respect to the update rule in equation (5);
19 Return  $\hat{w}$ ;
```

57 4 Algorithm: 2D and 3D action-space

58 **Note**, here we are exploring (chain components) the action space therefore the step-size b_o represents
 59 the length of the bond vector between actions.

60 Given the following rotation operators,

$$\prod_2^\theta = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}$$

$$\prod_3^\theta = \begin{bmatrix} \cos \theta + u_x^2 (1 - \cos \theta) & u_x u_y (1 - \cos \theta) - u_z \sin \theta & u_x u_z (1 - \cos \theta) + u_y \sin \theta \\ u_y u_x (1 - \cos \theta) + u_z \sin \theta & \cos \theta + u_y^2 (1 - \cos \theta) & u_y u_z (1 - \cos \theta) - u_x \sin \theta \\ u_z u_x (1 - \cos \theta) - u_y \sin \theta & u_z u_y (1 - \cos \theta) + u_x \sin \theta & \cos \theta + u_z^2 (1 - \cos \theta) \end{bmatrix}$$

62 around point $(0, 0, 0)$ and unit axis (u_x, u_y, u_z) . Update rule,

$$w_{t+1} \leftarrow w_t + \alpha [R_{t+1} + \gamma \max_a Q_{w_t}(S_{t+1}, a) - Q_{w_t}(S_t, A_t)] \nabla_{w_t} Q_{w_t}(S_t, A_t) \quad (5)$$

Algorithm 2: PolyRL: 2D

Input: $\mathcal{A}, \mathcal{S}, \prod_d^\bullet$, Intrinsic Rigidity κ^i , temperature \mathcal{E}_T , d , Step-size $b_o, \phi, \gamma, \rho, \epsilon, \alpha, \sigma, T$

Output: \hat{w}

```
1 Construct  $L_p^{(3)} \leftarrow \frac{\kappa^i}{\mathcal{E}_T}$  ;
2  $L_p^{(d)} \leftarrow \frac{2L_p^{(3)}}{(d-1)}$ ;
3  $\mu \leftarrow \cos^{-1}(e^{\frac{-b_o}{L_p^{(d)}}})$ ;
4 Let  $w_0 = 0$ ;
5 for each epoch do
6   if  $t == 0$  then
7     Sample  $A_0$  and  $S_0$  w.r.t  $\rho$ ;
8      $\theta \leftarrow$  uniform draw from  $(-\pi, \pi)$ ;
9   else if  $t == 1$  then
10     $A_1 \leftarrow (A_0(0) + b_o \cos(\theta), A_0(1) + b_o \sin(\theta))$ ;
11   else
12     Draw a sample  $\theta$  from  $\mathcal{N}(\mu, \sigma)$ ;
13      $\theta_t \leftarrow$  toss a coin and choose between  $\theta$  and  $-\theta$ ;
14      $H_{t-1} \leftarrow A_{t-1} - A_{t-2}$ ;
15      $A_t \leftarrow A_{t-1} +$  apply  $\prod_i^{\theta_t}$  on  $H_{t-1}$ ;
16   if  $A_t$  is valid then
17     Apply step function on action  $A_t$  and observe  $R_{t+1}$  and  $S_{t+1}$ ;
18   else
19     End the episode and re-start the chain;
20   Update  $\hat{w}$  with respect to the update rule in equation (5);
21 Return  $\hat{w}$ ;
```

Algorithm 3: PolyRL: Higher Dimension $d \geq 3$

Input: \mathcal{A}, \mathcal{S} , Persistence Length L_p , Action Space Dimension d , Step-size b_o, γ ,
Initial State-Action Distribution $\rho, \epsilon, \alpha, \sigma$, Episode Length T

Output: \hat{w}

```
1  $\mu \leftarrow \cos^{-1}(e^{\frac{-b_o}{L_p(d)}})$ ;
2 remainder  $\leftarrow d \bmod 3$ ;
3 if remainder==0 then
4   Subchain=d/3;
5   Group actions coordinates to “subchain” number of groups;
6 else
7   Subchain=(d+(3-remainder))/3;
8   Group actions to “subchain” number of groups; add total of “3-remainder” dummy dimensions
   such that each group has at most “1” dummy dimension;
9 Let  $w_0 = 0$ ;
10 for each epoch till  $T$  do
11   for  $i=1$ :Subchain do
12     if  $t == 0$  then
13       Sample  $A_0(i)$  and  $S_0(i)$  w.r.t  $\rho$ ;
14        $\varphi(i) \leftarrow$  uniform draw from  $(0, 2\pi)$ ;
15        $\theta(i) \leftarrow$  uniform draw from  $(-\pi/2, \pi/2)$ ;
16     else if  $t == 1$  then
17        $A_1(i) \leftarrow (A_0(i, 1) + b_o \cos(\theta(i)) \sin(\varphi(i)), A_0(i, 2) +$ 
        $b_o \sin(\theta(i)) \sin(\varphi(i)), A_0(i, 3)) + b_o \cos(\varphi(i))$ ;
18     else
19       Draw a sample  $\theta(i)$  from  $\mathcal{N}(\mu, \sigma)$ ;
20        $\theta_t(i) \leftarrow$  toss a coin and choose between  $\theta(i)$  and  $-\theta(i)$ ;
21        $\varphi_t(i) \leftarrow$  uniform draw from  $(0, 2\pi)$ ;
22        $A_t(i) \leftarrow (A_{t-1}(i, 1) + b_o \cos(\theta_t(i)) \sin(\varphi_t(i)), A_{t-1}(i, 2) +$ 
        $b_o \sin(\theta_t(i)) \sin(\varphi_t(i)), A_{t-1}(i, 3)) + b_o \cos(\varphi_t(i))$ ;
23   if  $A_t$  is valid then
24     Apply step function on action  $A_t$  and observe  $R_{t+1}$  and  $S_{t+1}$ ;
25   else
26     End the episode and re-start the chain;
27   Update  $\hat{w}$  with respect to the update rule in equation;
28   (5) Return  $\hat{w}$ ;
```
