
World Model as a Graph: Learning Latent Landmarks for Planning Supplementary Materials

Lunjun Zhang^{1,2} Ge Yang³ Bradly Stadie⁴

1. Greedy Latent Sparsification

Algorithm 2 Greedy Latent Sparsification (GLS) for Latent Cluster Training

Given: Replay Buffer B , Encoder f_E .

Initialize: LatentEmbeds = $\{\}$. ▷ Set of embeddings selected.

```
1: Sample  $K$  achieved goals from  $B$ .
2: Sample  $k \sim \{0 \dots K - 1\}$ .
3:  $\text{dist} = [\|f_E(g_1) - f_E(g_k)\|_2^2, \dots, \|f_E(g_K) - f_E(g_k)\|_2^2]$ 
4: for  $i = 1$  to  $M$  do ▷ Sub-sampling
5:    $k \leftarrow \arg \max \text{dist}[k]$ 
6:   Add  $f_E(g_k)$  to LatentEmbeds.
7:    $\text{NEWdist} = [\|f_E(g_1) - f_E(g_k)\|_2^2, \dots, \|f_E(g_K) - f_E(g_k)\|_2^2]$ 
8:    $\text{dist} = \text{ElementwiseMin}(\text{dist}, \text{NEWdist})$ 
9: end for
10: Optimize equation 5 on LatentEmbeds.
```

The Greedy Latent Sparsification (GLS) algorithm subsamples a large batch by sparsification. GLS first randomly selects a latent embedding from the batch, and then greedily chooses the next embedding that is furthest away from already selected embeddings. After collecting some *warm-up trajectories* before planning starts (see Table 1 below) during training, we first use GLS to initialize the latent centroids, and then continue to use it to sample the batches used to train the latent clusters. GLS is strongly inspired by (Arthur & Vassilvitskii, 2007), and this type of approach is known to improve clustering.

2. Graph Search with Soft Relaxations

In this paper, we employ a soft version of Floyd algorithm, which we find to empirically work well. Rather than simply using the min operation to do relaxation, the soft value iteration procedure uses a *soft* min operation when doing an update (note that, since we negated the distances to be negative in the weight matrix of the graph, the operations we use are actually max and softmax). The reason is that neural distances can be inconsistent and inaccurate at times, and using a soft operation makes the whole procedure more

robust. More concretely, we repeat the following update on the weight matrix for S steps with temperature β :

$$w_{i,j} \leftarrow \sum_{k=1}^{N+1} \frac{\exp \frac{1}{\beta}(w_{i,k} + w_{k,j})}{\sum_{k'=1}^{N+1} \exp \frac{1}{\beta}(w_{i,k'} + w_{k',j})} (w_{i,k} + w_{k,j}) \quad (1)$$

Following the practice in (Eysenbach et al., 2019; Huang et al., 2019), we do the following initialization to the distance matrix: for entries smaller than the negative of d_{max} , we penalize the entry by adding $-\infty$ to it (in this paper, we use -10^6 as the $-\infty$ value). The essential idea is that we only trust a neural estimate when it is *local*, and we rely on graph search to solve for *global*, longer-horizon distances. The $-\infty$ penalty effectively masks out those entries with large negative values in the softmax operation above. If we replace softmax with a hard max, we recover the original update in Floyd algorithm; we can interpolate between a hard Floyd and a soft Floyd by tuning the temperature β .

3. Overall Training Procedure

Here we provide an overall training procedure for L^3P in **Algorithm 3**. Given an environment env and a training goal distribution $p(g)$, we initialize a replay buffer B and the following **trainable modules**: policy π , distance function D , value function V , encoder f_E and decoder f_D , latent centroids $\{c_1 \dots c_N\}$.

Every K_{env} episodes of sampling, we take gradient steps for the above modules. The ratio between the number of environment steps and the number of gradient steps is a hyper-parameter.

4. Implementation Details

- We find that having a centralized replay for all parallel workers is significantly more sample efficient than having separate replays for each worker and simply averaging the gradients across workers.
- For Ant-Maze environment, we do grad norm clipping

¹University of Toronto ²Vector Institute ³MIT ⁴Toyota Technological Institute at Chicago. Correspondence to: Lunjun Zhang <lunjun@cs.toronto.edu>.

