

## A APPENDIX

### A.1 INFORMATION BASED MEASURE DERIVATION

Given that the choice of  $p_z$  for sampling  $z_i$ s is  $p(z) \sim \text{Uniform}(1, n_z)$  a categorical uniform distribution where the cardinality of  $|Z| = n_z$  we can write the mutual information based measure as follows.

$$\begin{aligned}
I(s; z) &= I(z_i; s_j) \\
I(z; s) &= \sum_{s_j, z_i} p(s_j, z_i) \log \frac{p(s_j, z_i)}{p(z_i)p(s_j)} \\
I(z; s) &= \sum_{s_j, z_i} p(s_j, z_i) \log \frac{p(z_i|s_j)}{p(z_i)} \\
I(z; s) &= \sum_{s_j, z_i} p(s_j|z_i) \cdot p(z_i) \log \frac{p(z_i|s_j)}{p(z_i)} \\
I(z; s) &= \sum_{s_j, z_i} p(s_j|z_i) \cdot n_z^{-1} \log \frac{p(z_i|s_j)}{n_z^{-1}} \\
I(z; s) &= \sum_{z_i} n_z^{-1} \cdot \sum_{s_j, z_i} d_{\pi_i}(s) \log \frac{d_{\pi_i}(s)}{\sum_k d_{\pi_k}(s) \cdot n_z^{-1}} \\
I(z; s) &= \sum_{z_i} n_z^{-1} \cdot KL(d_{\pi_{z_i}}(\cdot) || d_{\pi_{z_{all}}}(\cdot)) \\
I(z; s) &= E_{z_i \sim p(z)} [KL(d_{\pi_{z_i}}(\cdot) || d_{\pi_{z_{all}}}(\cdot))]
\end{aligned} \tag{13}$$

Here we can write  $p(s_j|z_i) = d_{\pi_i}(s_j)$  as it is the steady state probability induced by the policy corresponding to a certain  $z_i \in Z$ . Furthermore we replace  $p(z_i|s_j)$  in 13 with the following.

$$\begin{aligned}
p(z_i|s_j) &= \frac{p(s_j|z_i) \cdot p(z_i)}{\sum_k p(s_j|k) \cdot p(k)} \\
p(z_i|s_j) &= \frac{d_{\pi_i} \cdot p(z_i)}{\sum_k d_{\pi_k}(s_j) \cdot p(k)} \\
p(z_i|s_j) &= \frac{d_{\pi_i}(s_j) \cdot n_z^{-1}}{\sum_k d_{\pi_k}(s_j) \cdot n_z^{-1}} \\
p(z_i|s_j) &= \frac{d_{\pi_i}(s_j)}{\sum_k d_{\pi_k}(s_j)}
\end{aligned} \tag{14}$$

Since KL divergence is a convex function we can see that the information measure for the  $i$ th policy is convex on  $d_{\pi_i}$  induced by that policy  $\pi_i$ .

### A.2 INFORMATION BASED MEASURE'S DEVIATION FROM IDEAL MEASURE

$$\begin{aligned}
\delta^{\text{info}} &= E_{z_i \sim p(z)} [KL(d_{\pi_{z_i}}(\cdot) || d_{\pi_{z_i-1}}(\cdot))] - E_{z_i \sim p(z)} [KL(d_{\pi_{z_i}}(\cdot) || d_{\pi_{z_{all}}}(\cdot))] \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s) \log \frac{d_{\pi_{z_i}}(s)}{d_{\pi_{z_i-1}}(s)}] - E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s) \log \frac{d_{\pi_{z_i}}(s)}{d_{\pi_{z_{all}}}(s)}] \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s_j) (\log \frac{d_{\pi_{z_i}}(s_j)}{\sum_{k \neq i} d_{\pi_k}(s) \cdot (n-1)^{-1}} - \log \frac{d_{\pi_{z_i}}(s_j)}{\sum_k d_{\pi_k}(s_j) \cdot n^{-1}})] \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s_j) \log (\frac{\sum_k d_{\pi_k}(s) \cdot n^{-1}}{\sum_{k \neq i} d_{\pi_k}(s) \cdot (n-1)^{-1}})]
\end{aligned} \tag{15}$$

### A.3 PROPOSED METHOD'S DEVIATION FROM IDEAL MEASURE

$$\begin{aligned}
\delta^{\text{ours}} &= E_{z_i \sim p(z)} [KL(d_{\pi_{z_i}}(\cdot) || d_{\pi_{z_i-1}}(\cdot))] - E_{z_i \sim p(z)} [KL(\text{Uniform}(\cdot) || d_{\pi_{z_i-1}}(\cdot))] \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s) \log \frac{d_{\pi_{z_i}}(s)}{d_{\pi_{z_i-1}}(s)}] - E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s) \log \frac{n_s^{-1}}{d_{\pi_{z_i-1}}(s)}] \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s_j) (\log \frac{d_{\pi_{z_i}}(s_j)}{\sum_{k \neq i} d_{\pi_k}(s) \cdot (n-1)^{-1}} - \log \frac{n_s^{-1}}{\sum_{k \neq i} d_{\pi_k}(s) \cdot (n-1)^{-1}})] \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s_j) (\log(d_{\pi_{z_i}}(s_j)) + \log(n_s))] \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} d_{\pi_i}(s_j) \log(d_{\pi_{z_i}}(s_j))] + K \\
\delta^{\text{info}} &= E_{z_i \sim p(z)} [-H(d_{\pi_{z_i}}(\cdot))] + K
\end{aligned} \tag{16}$$

### A.4 COMPARISON OF OUR DEVIATION VS INFORMATION BASED MEASURE'S DEVIATION

Since log is a monotonic function we can lower bound the information based measure's deviation as follows

$$\begin{aligned}
\delta^{\text{info}} &= E_{z_i \sim p(z)} [\sum_{s_j} (d_{\pi_i}(s_j) \log(\frac{\sum_{k \neq i} d_{\pi_k}(s) \cdot (n-1)_z^{-1}}{\sum_k d_{\pi_k}(s) \cdot n_z^{-1}}))] \\
\delta^{\text{info}} &\leq E_{z_i \sim p(z)} [\sum_{s_j} (d_{\pi_i}(s_j) \log(\frac{\sum_{k \neq i} d_{\pi_k}(s) \cdot (n-1)_z^{-1}}{d_{\pi_i}(s) \cdot n_z^{-1}}))] \\
\delta^{\text{info}} &\leq E_{z_i \sim p(z)} [\sum_{s_j} (d_{\pi_i}(s_j) (\log(\frac{\sum_{k \neq i} d_{\pi_k}(s) \cdot (n-1)_z^{-1}}{d_{\pi_i}(s)}) + \log(n_z))] \\
\delta^{\text{info}} &\leq E_{z_i \sim p(z)} [-KL(d_{\pi_{z_i}}(\cdot) || d_{\pi_{z_i-1}}(\cdot))] + \log(n_z) \\
\delta^{\text{info}} &\geq E_{z_i \sim p(z)} [KL(d_{\pi_{z_i}}(\cdot) || d_{\pi_{z_i-1}}(\cdot))] - C \\
\delta^{\text{info}} &\geq E_{z_i \sim p(z)} [H(d_{\pi_{z_i}}(\cdot), d_{\pi_{z_i-1}}(\cdot)) - H(d_{\pi_{z_i}}(\cdot))] - C
\end{aligned} \tag{17}$$

If we compared both  $\delta^{\text{info}}$  and  $\delta^{\text{ours}}$

$$\begin{aligned}
\delta^{\text{info}} - \delta^{\text{ours}} &= H(d_{\pi_{z_i}}(\cdot), d_{\pi_{z_i-1}}(\cdot)) - H(d_{\pi_{z_i}}(\cdot)) + H(d_{\pi_{z_i}}(\cdot)) + M \\
\delta^{\text{info}} - \delta^{\text{ours}} &= H(d_{\pi_{z_i}}(\cdot), d_{\pi_{z_i-1}}(\cdot)) + M
\end{aligned} \tag{18}$$

For a discrete case the cross entropy is a non negative quantity. Thus always  $\delta^{\text{ours}} \leq \delta^{\text{info}}$  making our measure close to the ideal measure than the information based measure.

#### A.5 VISUALIZATION OF THE DIVERSITY MEASURE LANDSCAPE FOR THE GRID WORLD PROBLEM

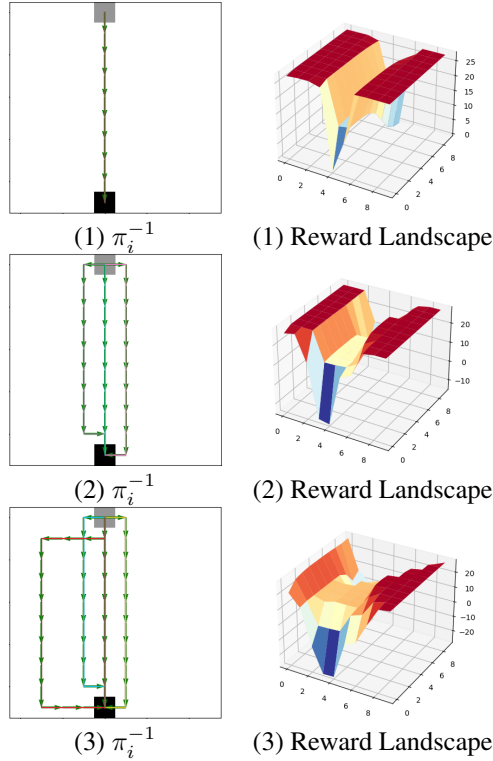


Figure 5: This figure illustrates learned diversity measure’s landscape corresponding policy of others in different occasion in the gridworld setting. For the figures in the right the x, y axis denote the grid and the z axis denotes the value of the reward function. This figure illustrates on how visited states are penalized by a function independent of the policy  $\pi_i$  under the proposed diversity measure.

## A.6 COVERAGE: CONTINUAL CONTROL TASK

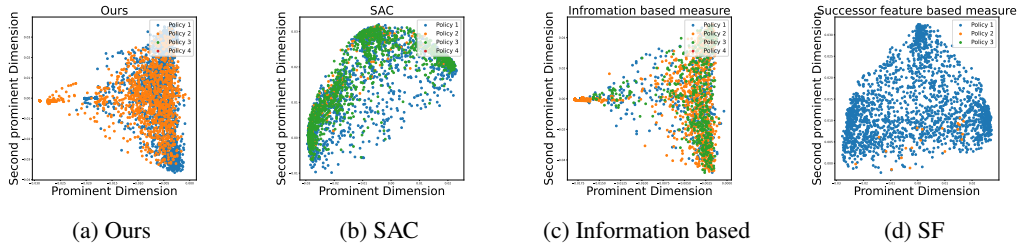


Figure 6: This figure illustrates the state space coverage by the policies generated by each agent in the ant environment. Dimensionality reduction was performed through PCA on the visited states. Here the x and y axes represents the two primary dimensions resulting from PCA. Our proposed method demonstrated better coverage in the areas on the right hand side of the space a compared to the lack of coverage in case of SAC. All three of the polices showing better performance than the standard SAC on test task while having certain level of coverage in this area leads us to hypothesize that the diverse polices are generated by the visitation of these states. Our proposed measure with a higher coverage was able perform better in average the than other methods both enhancing our hypothesis and the ability of our measure to induce better diversity.