
A New Representation of Successor Features for Transfer across Dissimilar Environments

Majid Abdolshah¹ Sunil Gupta¹ Santu Rana¹ Thommen George Karimpanal¹ Hung Le¹ Svetha Venkatesh¹

Abstract

Transfer in reinforcement learning is usually achieved through generalisation across tasks. Whilst many studies have investigated transferring knowledge when the reward function changes, they have assumed that the dynamics of the environments remain consistent. Many real-world RL problems require transfer among environments with different dynamics. To address this problem, we propose an approach based on successor features in which we model successor feature functions with Gaussian Processes permitting the source successor features to be treated as noisy measurements of the target successor feature function. Our theoretical analysis proves the convergence of this approach as well as the bounded error on modelling successor feature functions with Gaussian Processes in environments with both different dynamics and rewards. We demonstrate our method on benchmark datasets and show that it outperforms current baselines.

1. Introduction

Reinforcement learning (RL) is a computational approach that learns how to attain a complex goal by maximising rewards over time. Successful applications range from Atari games (Mnih et al. (2015)), to robotics (Zhang et al. (2017)), and self-driving cars (Liang et al. (2018)). However, this success is based on solving each task from scratch, and thus training these agents requires vast amounts of data.

Several solutions have been proposed to address this problem. Most works in transfer RL such as Progressive Neural Networks (Rusu et al. (2016)) and Inter-Task Mapping setups (Ammar & Taylor (2011); Gupta et al. (2017); Konidaris & Barto (2006); Yin & Pan (2017)) assume that the state-action space, or reward distribution space can be

¹Applied Artificial Intelligence Institute (A²I²), Deakin University, Geelong, Australia. Correspondence to: Majid Abdolshah <majid@deakin.edu.au>.

disentangled into independent sub-domains. However, learning interpretable, disentangled representations is challenging (Zhu et al. (2020)). The dynamics and the reward was decoupled for the first time in (Barreto et al. (2017)). Building upon an elegant formulation called the *successor function* (Dayan (1993)), the method allowed flexible transfer learning across tasks that differ in their reward structure. The underlying assumption was that the environmental dynamics remains unchanged and using a Generalised Policy Improvement (GPI) method, the optimal policy is determined. Such successor feature based methods have been shown to efficiently transfer knowledge across RL tasks (Barreto et al. (2018; 2019; 2020)). Other works that have built upon successor features include generalised policy updates on successor features (Barreto et al. (2020)), a universal type of successor feature based on the temporal difference method (Ma et al. (2020)), and Variational Universal Successor Features (Siriwardhana et al. (2019)) that perform target driven navigation. Option Keyboard (Barreto et al. (2019)) leveraged successor features to combine skills to define and manipulate options. This allows for change in reward, but a slight change of environment can deteriorate the performance of a new task. If however both the environmental dynamics and the reward differs across tasks, these methods are unable to handle this challenge.

In real-world problems when environment dynamics and rewards change across tasks, both these aspects need careful modelling. Failing to do so can lead to negative transfer of knowledge from the previously seen tasks. Thus RL methods using successor features need to be extended to handle the changes in environmental dynamics. Such work is limited. Zhang et al. (Zhang et al. (2017)) aim to address this problem by considering a linear relationship between the source and target successor features. This modelling is restrictive and may fall short in capturing the complexity of the changed environment dynamics. *Thus, the problem of designing a method using a successor feature based approach to transfer knowledge from source to target environment where the dynamics are dissimilar, is still open.*

Our new approach enables the efficient learning of novel successor features to cater to the new target environmental dynamics. This is done by using the distribution of the previous (source) task successor features as a prior for the new tar-

get task. We model both the target and source distributions through Gaussian Processes (GPs). The target distribution is modeled as a noisy version of the source distribution. This approach assumes that the source and target environments lie within some proximity to each other i.e. they are similar within some noisy envelope. However, this adjustable noisy envelope impacts the upper bounded error on the modelling of optimal policy in the target environment. The advantage of this approach is that the source observations provide a head-start for the learning process and additional explorations in the target will provide efficient convergence to the optimal policy. We use a GPI method to estimate the target action value function. We provide theoretical analysis and upper bounds (1) on the difference of action-value functions when the optimal policy derived from environment i is replaced by the optimal policy derived from environment j ; (2) on the estimation error of the action-value function of an optimal policy learned in a source environment when executed in the target environment; (3) on the difference of the optimal action-value function in the target environment and our GPI-derived action value function. We evaluate this approach in a variety of benchmark environments with different levels of complexity. Our key contributions are:

- A new method based on successor features and Gaussian Processes that enhances transfer from source to target tasks when the dynamics of the environment are dissimilar;
- A theoretical analysis for the new successor based method; and,
- An empirical comparison on diverse suit of RL benchmarks with different levels of complexity.

2. Background

2.1. Reinforcement Learning

We model the RL framework as a Markov Decision Process (MDP) described as $\langle \mathcal{S}, \mathcal{A}, p, R \rangle$, where \mathcal{S} is a finite state space, \mathcal{A} represents a finite action space, $p : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ the transition probabilities, and $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a bounded reward function. A discount factor $\gamma \in (0, 1]$ encodes the importance of future rewards with respect to the current state. The objective of RL is to find an optimal policy $\pi(a|s) : \mathcal{S} \rightarrow \mathcal{A}$ that maps the states to the actions such that it maximises the expected discounted reward. The action-value of policy π is defined as $Q^\pi(s_t, a_t) = \mathbb{E}^\pi[\sum_{k=0}^{\infty} \gamma^k R(s_{t+k}, a_{t+k}) | s_0 = s_t, a_0 = a_t]$, where s_t and a_t are the state of the agent at time step t , and the action that is taken in that state, respectively.

Q-Learning: Q-Learning is an off-policy RL approach that aims to learn the optimal action-value function. This function can be updated recursively as:

$$Q^\pi(s_t, a_t) = \mathbb{E}_{s_{t+1}} [R(s_t, a_t) + \gamma \max_{a \in \mathcal{A}} (Q^\pi(s_{t+1}, a))]. \quad (1)$$

After learning the action-value function, the optimal policy can be retrieved by selecting the best action at every state: $\pi^*(s) \in \operatorname{argmax}_{a \in \mathcal{A}} Q^*(s, a)$. In the next section we draw the connection between Q-Learning and successor features.

2.2. Successor Features

The successor feature representation allows decoupling of the dynamics of an MDP from its reward distributions. Barreto et al. (Barreto et al. (2017)) generalised the successor representations that was first formulated in Dayan et al. (Dayan (1993)) decomposing the action-value function into a set of features that encode the dynamics of the environment and a weight that acts as a task-specific reward mapper. This decomposition can be formulated as:

$$R(s, a) = \phi(s, a)^\top \mathbf{w}, \quad (2)$$

where $\phi(s, a) \in \mathbb{R}^D$ are the features of (s, a) that represent the dynamics of the environment and \mathbf{w} is the reward mapper of the environment. The two components can be learnt through supervised learning (Zhu et al. (2020)). Note that $\phi(\cdot, \cdot)$ can be any complex model such as a neural network. Barreto et al. (Barreto et al. (2017)) showed that this decomposition can be used in the construction of the action-value function. Let us assume a reward function as in Eq. (2), the action value function can be derived as:

$$\begin{aligned} Q^\pi(s, a) &= \mathbb{E}^\pi [R(s_{t+1}, a_{t+1}) + \gamma R(s_{t+2}, a_{t+2}) + \\ &\dots | s_t = s, a_t = a] = \mathbb{E}^\pi \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_{t+1}, a_{t+1}) \right. \\ &\left. | s_t = s, a_t = a \right] \mathbf{w} = \psi(s, a)^\top \mathbf{w}. \end{aligned} \quad (3)$$

where $\psi(s, a)^\top$ is the “Successor Feature (SF)” of (s, a) and summarises the dynamics induced by π . Eq. (3) satisfies the Bellman Equation and can be learnt through any conventional method. By treating the latent representation $\phi(\cdot, \cdot)$ as the immediate reward in the context of Q-Learning, the successor feature function can be written as:

$$\psi^\pi(s, a) = \phi(s_t, a_t) + \gamma \mathbb{E}^\pi [\psi^\pi(s_{t+1}, \pi(s_{t+1})) | s_t = s, a_t = a]. \quad (4)$$

The principal advantage of SFs is that when the knowledge of $\psi^\pi(s, a)$ is observed, one can compute a task-specific reward mapper based on the observations seen in the same environment with different reward function. Given $\psi^\pi(s, a)$, the updated reward mapper $\tilde{\mathbf{w}}$ for a new reward function can be approximated by solving a regression problem in tabular scenarios (Barreto et al. (2020)). This enables us to construct

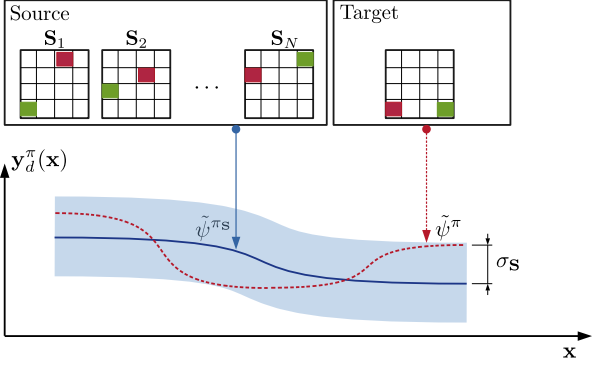


Figure 1. Our proposed approach uses GPs to model the source successor features functions ($\tilde{\psi}^{\pi_s}$) as noisy measurements for the target successor features functions ($\tilde{\psi}^{\pi}$).

the new action-value function $\tilde{Q}^{\pi}(s, a) = \psi^{\pi}(s, a)^{\top} \tilde{\mathbf{w}}$ by only observing few steps of the new reward function in a similar environment. In non-tabular problems, a deep neural network can be used to learn the successor feature functions of an environment by minimising the following loss:

$$L_{\psi}(\boldsymbol{\theta}) = \|\tilde{\phi}(s_t, a_t) + \gamma^t \tilde{\psi}^{\pi}(s_{t+1}, a_{t+1}; \boldsymbol{\theta}_{\psi}) - \tilde{\psi}^{\pi}(s_t, a_t; \boldsymbol{\theta}_{\psi})\|. \quad (5)$$

Similarly, by having the $\phi(\cdot, \cdot)$ function, the new task-specific reward mapper $\tilde{\mathbf{w}}$ is computed as:

$$\mathcal{L}_{\tilde{\mathbf{w}}}(\tilde{\mathbf{w}}) = \mathbb{E}_{\forall (s, a) \in \text{new task}} \left[(r(s, a) - \phi(s, a)^{\top} \tilde{\mathbf{w}})^2 \right]. \quad (6)$$

where $r(s, a)$ is the obtained reward.

3. Method

We now consider source and target environments with both *dissimilar dynamics* and *different reward functions*. The goal of transfer in RL in such problems is to learn an optimal policy for the target environment, by leveraging exterior source information and interior target information (Joy et al. (2019); Shilton et al. (2017)). We first define the set of source environments as:

$$\mathcal{M}^{\mathbf{S}} = \left\{ \mathcal{M}(\mathcal{S}, \mathcal{A}, p_1, R_1), \dots, \mathcal{M}(\mathcal{S}, \mathcal{A}, p_N, R_N) \right\},$$

where $\mathcal{M}(\cdot)$ represents an MDP induced by $\phi(\cdot, \cdot)$ as a feature function used in all environments. We denote the N source environments by $\mathbf{S}_1, \dots, \mathbf{S}_N$ and train an optimal policy in each environment to create a set of policies: $\Pi^{\mathbf{S}} = \{\pi^{\mathbf{S}_1}, \dots, \pi^{\mathbf{S}_N}\}$. By executing these N optimal policies in their corresponding source environments, N distinct successor feature functions ψ^{π_i} , $i = \{1, \dots, N\}$ are computed using the loss function Eq. (5). We define $\psi^{\pi_i} = \psi_1^{\pi_i}, \dots, \psi_D^{\pi_i}$ and $\psi_d^{\pi_i}$ is the d -th dimension. Source

Algorithm 1 Extracting Successor Feature Functions of Source Environments.

- 1: **Input:**
- 2: Source observation and learned policies $\mathcal{D}^{\mathbf{S}_{i=1 \dots N}} = \{\cdot\}, \Pi^{\mathbf{S}}$.
- 3: Discount factor γ , exploration rate ε_e .
- 4: Feature function $\phi(\cdot, \cdot) \in \mathbb{R}^D$.
- 5: M_{\max} maximum number of episodes.e
- 6: **Output:** Successor feature functions in source environments $\tilde{\psi}^{\pi_1}, \dots, \tilde{\psi}^{\pi_N}$.
- 7: Initialise $\tilde{\psi}^{\pi_{1 \dots N}} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^D$.
- 8: **for** $i \in 1, \dots, N$ **do**
- 9: **while** M_{\max} **do**
- 10: Sample the initial state randomly $s \in \mathcal{S}$.
- 11: **while** $t \in$ steps not terminated **do**
- 12: **if** $\varepsilon -$ greedy **then**
- 13: $a_t = \text{Uniform}(\mathcal{A})$. //random action
- 14: **else**
- 15: $a_t = \pi^{\mathbf{S}_i}(s_t)$.
- 16: **end if**
- 17: Execute a_t , observe s_{t+1} .
- 18: $a_{t+1} = \pi^{\mathbf{S}_i}(s_{t+1})$.
- 19: Minimise the loss $L_{\psi^{\pi_i}}(\boldsymbol{\theta})$ //Eq. (5)
- 20: Store the transition $[s_t, a_t, \tilde{\psi}^{\pi_i}(s_t, a_t; \boldsymbol{\theta}_{\psi^{\pi_i}})]$ in $\mathcal{D}^{\mathbf{S}_i}$.
- 21: Perform gradient descent w.r.t. $\boldsymbol{\theta}_{\psi^{\pi_i}}$.
- 22: **end while**
- 23: **end while**
- 24: **end for**

successor feature functions can be learnt using neural networks with parameters $\boldsymbol{\theta}_{\psi}$ that can be updated by gradient descent. Algorithm 1 shows how the successor feature functions are computed.

Using state-action pairs as observations, we construct a set of successor feature function samples as: $\mathcal{D}^{\mathbf{S}_i} = \left\{ (\mathbf{x}_1^{\mathbf{S}_i}, \tilde{\psi}^{\pi_i}(\mathbf{x}_1^{\mathbf{S}_i})), (\mathbf{x}_2^{\mathbf{S}_i}, \tilde{\psi}^{\pi_i}(\mathbf{x}_2^{\mathbf{S}_i})), \dots, (\mathbf{x}_n^{\mathbf{S}_i}, \tilde{\psi}^{\pi_i}(\mathbf{x}_n^{\mathbf{S}_i})) \right\}$, where $\mathbf{x}_t^{\mathbf{S}_i} = (s_t, a_t)$ is a tuple of state-action in \mathbf{S}_i at time t following policy π_i and $\tilde{\psi}^{\pi_i}(\mathbf{x}_t^{\mathbf{S}_i})$ is the successor feature for $\mathbf{x}_t^{\mathbf{S}_i}$. If observations from target environment exist, a set of successor feature function samples in the target environment can be constructed as $\mathcal{D}^{\mathcal{T}} = \left\{ (\mathbf{x}_1^{\mathcal{T}}, \tilde{\psi}^{\pi}(\mathbf{x}_1^{\mathcal{T}})), (\mathbf{x}_2^{\mathcal{T}}, \tilde{\psi}^{\pi}(\mathbf{x}_2^{\mathcal{T}})), \dots, (\mathbf{x}_m^{\mathcal{T}}, \tilde{\psi}^{\pi}(\mathbf{x}_m^{\mathcal{T}})) \right\}$, where $\mathbf{x}_t^{\mathcal{T}} = (s_t, a_t)$ records a tuple of state-action visited at time step t following policy π in the target environment \mathcal{T} , and $\tilde{\psi}^{\pi}(\mathbf{x}_t^{\mathcal{T}})$ is an approximation of successor feature of $\mathbf{x}_t^{\mathcal{T}}$ in the target environment.

We assume that the target successor feature function $\tilde{\psi}_d^{\pi}$, following policy π has a measurement noise of $\epsilon_d \sim \mathcal{N}(0, \sigma^2)$, i.e.:

$$y_d^{\pi}(\mathbf{x}^{\mathcal{T}}) = \tilde{\psi}_d^{\pi}(\mathbf{x}^{\mathcal{T}}) + \epsilon_d, \quad \forall d \in \{1, \dots, D\},$$

where $\mathbf{y}_d^\pi(\mathbf{x}^\mathcal{T})$ represents the noisy value of successor features in the target environment. Likewise, source successor features are assumed to have a measurement noise of $\epsilon_d \sim \mathcal{N}(0, \sigma^2)$, i.e.:

$$\mathbf{y}_d^{\pi^i}(\mathbf{x}^{\mathbf{S}^i}) = \tilde{\psi}_d^{\pi^i}(\mathbf{x}^{\mathbf{S}^i}) + \epsilon_d, \forall d \in \{1, \dots, D\}, i = \{1, \dots, N\}.$$

We model the samples of successor features from the source environment as noisy variants of successor features in the target environment. Under this model, the target successor feature for observation $\mathbf{x}^\mathcal{T}$ is defined as:

$$\tilde{\psi}_d^\pi(\mathbf{x}^\mathcal{T}) = \tilde{\psi}_d^{\pi^i}(\mathbf{x}^\mathcal{T}) + \epsilon_d^{\mathbf{S}^i}, \forall d \in \{1, \dots, D\},$$

where $\epsilon_d^{\mathbf{S}^i}$ is the modelling noise of target successor feature functions. We assume the modeling noise to be Gaussian distributed with variance $\sigma_{\mathbf{S}}^2$ as $\epsilon_d^{\mathbf{S}^i} = \epsilon_d^{\mathbf{S}} \sim \mathcal{N}(0, \sigma_{\mathbf{S}}^2)$, $i = \{1, \dots, N\}$. Intuitively, this allows the use of the source samples as a noisy version for the target successor feature values. The value of the ‘‘modeling noise’’ variance $\sigma_{\mathbf{S}}^2$ depends on the difference between source and target (Shilton et al. (2017)) environments. Using the successor feature samples from both source and target environments, we model the target successor feature function using a Gaussian Process. The overall idea is illustrated in Figure 1. Without loss of generality, we use $\mathcal{GP}(0, k(\cdot, \cdot))$, i.e. a GP with a zero mean function, a symmetric positive-definite covariance function $k(\mathbf{x}, \mathbf{x}') : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$ a.k.a. kernel, and we also assume $k(\mathbf{x}, \mathbf{x}) = 1, \forall \mathbf{x} \in \mathbb{X}$. A covariance matrix \mathbf{K} based on the combined source and target samples of successor features can be written as:

$$\mathbf{K} = \begin{bmatrix} k(\mathbf{x}_1^{\mathbf{S}}, \mathbf{x}_1^{\mathbf{S}}) & \dots & k(\mathbf{x}_1^{\mathbf{S}}, \mathbf{x}_n^{\mathbf{S}}) & k(\mathbf{x}_1^{\mathbf{S}}, \mathbf{x}_1^{\mathcal{T}}) & \dots & k(\mathbf{x}_1^{\mathbf{S}}, \mathbf{x}_m^{\mathcal{T}}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_n^{\mathbf{S}}, \mathbf{x}_1^{\mathbf{S}}) & \dots & k(\mathbf{x}_n^{\mathbf{S}}, \mathbf{x}_n^{\mathbf{S}}) & k(\mathbf{x}_n^{\mathbf{S}}, \mathbf{x}_1^{\mathcal{T}}) & \dots & k(\mathbf{x}_n^{\mathbf{S}}, \mathbf{x}_m^{\mathcal{T}}) \\ k(\mathbf{x}_1^{\mathcal{T}}, \mathbf{x}_1^{\mathbf{S}}) & \dots & k(\mathbf{x}_1^{\mathcal{T}}, \mathbf{x}_n^{\mathbf{S}}) & k(\mathbf{x}_1^{\mathcal{T}}, \mathbf{x}_1^{\mathcal{T}}) & \dots & k(\mathbf{x}_1^{\mathcal{T}}, \mathbf{x}_m^{\mathcal{T}}) \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_m^{\mathcal{T}}, \mathbf{x}_1^{\mathbf{S}}) & \dots & k(\mathbf{x}_m^{\mathcal{T}}, \mathbf{x}_n^{\mathbf{S}}) & k(\mathbf{x}_m^{\mathcal{T}}, \mathbf{x}_1^{\mathcal{T}}) & \dots & k(\mathbf{x}_m^{\mathcal{T}}, \mathbf{x}_m^{\mathcal{T}}) \end{bmatrix}, \quad (7)$$

where $k(\mathbf{x}_i^{\mathbf{S}}, \mathbf{x}_j^{\mathbf{S}})$, $i, j = \{1, \dots, n\}$ is the self-covariance among source observations of $\mathbf{S} \in \{\mathbf{S}_1, \dots, \mathbf{S}_N\}$ and $k(\mathbf{x}_i^{\mathbf{S}}, \mathbf{x}_j^{\mathcal{T}})$, $i = \{1, \dots, n\}, j = \{1, \dots, m\}$ denotes the covariance between source and target observations. After incorporating the source, target measurement noise with the modeling noise, the covariance matrix can be written as:

$$\mathbf{K}_* = \mathbf{K} + \begin{bmatrix} (\sigma_{\mathbf{S}}^2 + \sigma^2)\mathbf{I}_{n \times n} & \mathbf{0} \\ \mathbf{0} & \sigma^2\mathbf{I}_{m \times m} \end{bmatrix}. \quad (8)$$

Intuitively, a higher value of $\sigma_{\mathbf{S}}^2$ implies higher uncertainty about similarity between source and target environments. Having defined the required components, using the property of GP (Rasmussen (2003)), the predictive mean and variance for a new target observation $\mathbf{x}^\mathcal{T} = (s, a)$ is derived as:

$$\mu_{m,d}(\mathbf{x}^\mathcal{T}) = \mathbf{k}^\mathbf{T} \mathbf{K}_*^{-1} \mathbf{y}_d \quad (9)$$

$$\sigma_{m,d}^2(\mathbf{x}^\mathcal{T}) = k(\mathbf{x}^\mathcal{T}, \mathbf{x}^\mathcal{T}) - \mathbf{k}^\mathbf{T} \mathbf{K}_*^{-1} \mathbf{k} \quad (10)$$

where \mathbf{K}_* is the kernel matrix as defined in Eq. (8) and $\mathbf{k} = [k(\mathbf{x}_i, \mathbf{x}^\mathcal{T})]$, $\forall \mathbf{x}_i \in \mathcal{D}^{\mathbf{S}} \cup \mathcal{D}^{\mathcal{T}}$. We use the posterior mean as in Eq. (9) as the predicted value of the d -th successor feature dimension for the target observation as $\tilde{\psi}_d^\pi(\mathbf{x}^\mathcal{T}) = \mu_{m,d}(\mathbf{x}^\mathcal{T})$.

Using GPI (Barreto et al. (2017)), we identify the optimal policy by selecting the best action of the best policy as:

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} \max_{\pi \in \Pi^{\mathbf{S}}} \tilde{Q}^\pi(s, a). \quad (11)$$

We note that $\tilde{Q}^\pi(\mathbf{x}^\mathcal{T}) \approx \tilde{\psi}^\pi(\mathbf{x}^\mathcal{T})^\mathbf{T} \tilde{\mathbf{w}}$, where $\tilde{\mathbf{w}}$ is obtained by minimising the loss function in Eq. (6) as the agent interacts with the target environment. We term our method *Successor Features for Dissimilar Environments (SFDE)* and is detailed in Algorithm 2.

3.1. Theoretical Analysis

This section answers the key question of ‘‘what are the effects of relaxing the assumption of similarity among source and target environments on the convergence and transfer via successor features?’’. We prove (1) an upper bound on the difference of action-value functions when the optimal policy derived from environment i is replaced by the optimal policy derived from environment j (Theorem 1); (2) an upper bound on the estimation error of the action-value function of an optimal policy learned in \mathbf{S}_j when executed in target environment \mathcal{T} (Lemma 1); (3) Using (1) and (2) an upper bound on the difference of the optimal action-value function in the target environment and our GPI-derived action value function (Theorem 2).

Theorem 1

Let \mathbf{S}_i and \mathbf{S}_j be two different source environments with dissimilar transition dynamics p_i and p_j respectively. Let $\delta_{ij} \triangleq \max_{s,a} |r_i(s, a) - r_j(s, a)|$, where $r_i(\cdot, \cdot)$ and $r_j(\cdot, \cdot)$ are the reward functions of environment \mathbf{S}_i and \mathbf{S}_j respectively. We denote π_i^* and π_j^* as optimal policies in \mathbf{S}_i and \mathbf{S}_j . It can be shown that the difference of their action-value functions is upper bounded as:

$$\begin{aligned} Q_i^{\pi_i^*}(s, a) - Q_i^{\pi_j^*}(s, a) &\leq \frac{2\delta_{ij}}{1-\gamma} \\ &+ \frac{\gamma \left\| \mathbf{P}_i(s, a) - \mathbf{P}_j(s, a) \right\|}{(1-\gamma)} \\ &\times \frac{\left(\left\| \mathbf{Q}_i^i - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_i^j \right\| \right)}{(1-\gamma)}, \quad (12) \end{aligned}$$

where $Q_i^{\pi_i^*}$ shows the action-value function in environment \mathbf{S}_i by following an optimal policy that is learned in the environment $\mathbf{S}_k \in \{\mathbf{S}_1, \dots, \mathbf{S}_N\}$. We also define $\mathbf{P}_i(s, a) = [p_i(s'|s, a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{P}_j(s, a) = [p_j(s'|s, a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_i^i = [\max_{b \in \mathcal{A}} Q_i^{\pi_i^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_j^j = [\max_{b \in \mathcal{A}} Q_j^{\pi_j^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_i^j = [\max_{b \in \mathcal{A}} Q_i^{\pi_j^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_j^i = [\max_{b \in \mathcal{A}} Q_j^{\pi_i^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, γ as the discount factor, and $\|\cdot\|$ to be 2-norm (Euclidean norm).

Proof: We provide a sketch of the proof which involves two key steps. The left side of the inequality (15) can be rewritten as:

$$\begin{aligned} Q_i^{\pi_i^*}(s, a) - Q_i^{\pi_j^*}(s, a) &= Q_i^{\pi_i^*}(s, a) - Q_j^{\pi_j^*}(s, a) \\ &\quad + Q_j^{\pi_j^*}(s, a) - Q_i^{\pi_j^*}(s, a) \\ &\leq \underbrace{|Q_i^{\pi_i^*}(s, a) - Q_j^{\pi_j^*}(s, a)|}_{(I)} \\ &\quad + \underbrace{|Q_j^{\pi_j^*}(s, a) - Q_i^{\pi_j^*}(s, a)|}_{(II)}. \end{aligned}$$

We can prove that (I) $\leq \frac{\delta_{ij}}{1-\gamma} + \gamma \|\mathbf{P}_i - \mathbf{P}_j\| \times \|\mathbf{Q}_i^i - \mathbf{Q}_j^j\| / (1-\gamma)$ and (II) $\leq \frac{\delta_{ij}}{1-\gamma} + \gamma \|\mathbf{P}_i - \mathbf{P}_j\| \times \|\mathbf{Q}_j^j - \mathbf{Q}_i^i\| / (1-\gamma)$, leading to the upper bound. Detailed proof is available in supplementary material.

Theorem 1 uses δ_{ij} as a metric of maximum immediate reward dissimilarities in the environments \mathbf{S}_i and \mathbf{S}_j . Clearly, the higher δ_{ij} , the less similar \mathbf{S}_i and \mathbf{S}_j are, and the upper bound will be looser accordingly. This upper bound also depends on the value of $\|\mathbf{P}_i - \mathbf{P}_j\|$ that captures the difference in dynamics of \mathbf{S}_i and \mathbf{S}_j - larger the value, looser the upper bound. However, $\|\mathbf{Q}_i^i - \mathbf{Q}_j^j\| + \|\mathbf{Q}_j^j - \mathbf{Q}_i^i\|$ incorporates the difference of action-value functions that are related to both the dynamics and the future discounted reward in the two environments. Hence, a larger value of this term implies that \mathbf{S}_i and \mathbf{S}_j are expected to produce different sum of discounted future reward by following their corresponding policies. Note that if $\mathbf{S}_i = \mathbf{S}_j$ (in terms of both dynamics and reward), the upper bound will vanish to 0 as the two environments are identical. Clearly, our bound is an extension of the bound in (Barreto et al. (2018)) for environments with dissimilar dynamics. In the special case of identical environments i.e. when $\|\mathbf{P}_i - \mathbf{P}_j\| = 0$, the two bounds become the same.

We now prove an upper bound on the estimation error of the action-value function of an optimal policy learned in \mathbf{S}_j when executed in target environment \mathcal{T} .

Algorithm 2 Successor Features for Dissimilar Environments.

```

1: Input:
2: Source environments  $\mathcal{D}^{\mathbf{S}_1, \dots, \mathbf{S}_N}$  and target observations  $\mathcal{D}^{\mathcal{T}} = \{\}$ .
3: Set the amount of noises for source and target  $\sigma_{\mathbf{S}}^2, \sigma^2$ .
4: SFs of source environments  $\tilde{\psi}^{\pi_1}, \dots, \tilde{\psi}^{\pi_N}$ .
5: Feature function  $\phi(\cdot, \cdot) \in \mathbb{R}^D$ .
6: Initialise reward mapper weight for target environment  $\tilde{\mathbf{w}}$ .
7: Output: Optimal policy  $\pi^*(s)$  for the target environment.
8: while  $t \in$  steps not terminated do
9:   for  $\forall \mathbf{S} \in \{\mathbf{S}_1, \dots, \mathbf{S}_N\}$  do
10:    for  $a' \in \mathcal{A}$  do
11:      $\mathbf{x}_t^{\mathcal{T}} = (s_t, a')$ 
12:      $\mathcal{GP}_d^{\mathbf{S}}(\{\mathcal{D}^{\mathbf{S}}, \mathcal{D}^{\mathcal{T}}\})$ ,  $\forall d = \{1, \dots, D\}$  //Fit the GPs with source and target data
13:
14:     Calculate  $\mathbf{K}_*, \mathbf{k}$ , //(Eq. 8)
15:      $\tilde{\psi}_d^{\mathbf{S}}(\mathbf{x}_t^{\mathcal{T}}) = \mu_{m,d}^{\mathbf{S}}(\mathbf{x}_t^{\mathcal{T}})$ ,  $\forall d = \{1, \dots, D\}$ , //(Eq. 9)
16:   end for
17: end for
18: Reconstruct the action-value functions of all source environments given the target observations  $\tilde{Q}^{\pi_i}(s_t, a') = \tilde{\psi}^{\pi_i}(s_t, a') \tilde{\mathbf{w}}$ ,  $\forall a' \in \mathcal{A}, \forall i = \{1, \dots, N\}$ .
19:  $\pi^*(s) = \operatorname{argmax}_{a' \in \mathcal{A}} \max_{\pi \in \Pi^{\mathbf{S}}} \tilde{Q}^{\pi}(\mathbf{x}_t^{\mathcal{T}})$ . //Performing GPI
20: Add the new target observation to  $\mathcal{D}^{\mathcal{T}}$ .
21: Update the estimation of  $\tilde{\mathbf{w}}$  based on Eq. (6).
22: end while

```

Lemma 1

Let π_1^*, \dots, π_N^* be N optimal policies for $\mathbf{S}_1, \dots, \mathbf{S}_N$ respectively and $\tilde{Q}_{\mathcal{T}}^{\pi_j^*} = (\tilde{\psi}^{\pi_j^*})^{\top} \tilde{\mathbf{w}}_{\mathcal{T}}$ denote the action-value function of an optimal policy learned in \mathbf{S}_j and executed in target environment \mathcal{T} . Let $\tilde{\psi}^{\pi_j^*}$ denote the estimated successor feature function from the combined source and target observations from \mathbf{S}_j and \mathcal{T} as defined in Eq. (9), and $\tilde{\mathbf{w}}_{\mathcal{T}}$ is the estimated reward mapper for environment \mathcal{T} by using loss function in Eq. (6). It can be shown that the difference of the true action-value function and the estimated one through successor feature functions and reward mapper, is bounded as:

$$\Pr\left(\left|Q_{\mathcal{T}}^{\pi_j^*}(s, a) - \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(s, a)\right| \leq \varepsilon(m) \forall s, a\right) \geq 1 - \delta,$$

where $\varepsilon(m) = \sqrt{2 \log(\frac{|\mathbb{X}| u_m}{\delta})} \sigma_{m,d}(\mathbf{x})$, $\mathbf{x} \in \mathbb{X}$, $\delta \in (0, 1)$, $u_m = \frac{\pi^2 m^2}{6}$, m being the number of observations in environment \mathcal{T} , and $\mathbf{x} = (s, a)$. $\sigma_{m,d}(\mathbf{x})$ is the square root

of posterior variance as defined in Eq. (10).

Proof: Proof is available in the supplementary material.

Lemma 1 ensures that reconstructing the action-value function on a new target environment \mathcal{T} by using $\tilde{\psi}^{\pi_j^*}$ and $\tilde{\mathbf{w}}_{\mathcal{T}}$ can be achieved with a bounded error with high probability. Essentially, the key term in $\varepsilon(m)$ is $\sigma_{m,d}^2(\mathbf{x})$ that is computed by using \mathbf{K}_* (see Eq. (8)), which itself incorporates the modeling noise variance $\sigma_{\mathbf{S}}^2$. Hence, by increasing $\sigma_{\mathbf{S}}^2$, $\sigma_{m,d}^2(\mathbf{x})$ will be higher and accordingly the upper bound will become looser. We note that due to the consistency of GPs, as $m \rightarrow \infty$, the uncertainty of predictions tends to 0 ($\sigma_{m,d}^2(\mathbf{x}) \rightarrow 0$) and thus $\varepsilon(m) \rightarrow 0$.

We now present our final result that bounds the difference of the optimal action-value function in the target environment and our GPI-derived action value function.

Theorem 2

Let $\mathbf{S}_{i=1\dots N}$ be N different source environments with dissimilar transition functions $p_{i=1\dots N}$. Let us denote the optimal policy π that is defined based on the GPI as:

$$\pi(s) \in \operatorname{argmax}_{a \in \mathcal{A}} \max_{j \in \{1 \dots N\}} \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(s, a), \quad (13)$$

where $\tilde{Q}_{\mathcal{T}}^{\pi_j^*} = (\tilde{\psi}^{\pi_j^*})^T \tilde{\mathbf{w}}_{\mathcal{T}}$ being the action-value function of an optimal policy learned in \mathbf{S}_j and executed in target environment \mathcal{T} , $\tilde{\psi}^{\pi_j^*}$ is the estimated successor feature from the combined source and target observations from \mathbf{S}_j and \mathcal{T} as defined in Eq. (9), and $\tilde{\mathbf{w}}_{\mathcal{T}}$ is the estimated reward mapper for target environment from Eq. (6). Considering Lemma 1 and Eq. (17), the difference of optimal action-value function in the target environment and our GPI-derived action value function is upper bounded as:

$$\begin{aligned} Q_{\mathcal{T}}^*(s, a) - \tilde{Q}_{\mathcal{T}}^{\pi \in \pi_j^*}(s, a) &\leq \frac{2\phi_{\max}}{1-\gamma} \left\| \tilde{\mathbf{w}}_{\mathcal{T}} - \mathbf{w}_j \right\| \\ &+ \frac{\gamma \left\| \mathbf{P}_{\mathcal{T}}(s, a) - \mathbf{P}_j(s, a) \right\|}{(1-\gamma)} \\ &\times \frac{\left(\left\| \mathbf{Q}_{\mathcal{T}}^* - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_{\mathcal{T}}^j \right\| \right)}{(1-\gamma)} \\ &+ \frac{2\varepsilon(m)}{(1-\gamma)}. \end{aligned} \quad (14)$$

where $\phi_{\max} = \max_{s,a} \|\phi(s, a)\|$. We also define $\mathbf{P}_{\mathcal{T}} = [p_{\mathcal{T}}(s'|s, a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{P}_j = [p_j(s'|s, a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_{\mathcal{T}}^* = [\max_{b \in \mathcal{A}} Q_{\mathcal{T}}^*(s', b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_j^j = [\max_{b \in \mathcal{A}} \tilde{Q}_j^{\pi_j^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_{\mathcal{T}}^j = [\max_{b \in \mathcal{A}} \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, and γ as the discount factor.

Proof: Proof is available in supplementary material.

In inequality (18), $\left\| \tilde{\mathbf{w}}_{\mathcal{T}} - \mathbf{w}_j \right\|$ encodes the dissimilarity of reward functions in target environment \mathcal{T} and \mathbf{S}_j , and $\varepsilon(m)$ holds the error of action-value reconstruction by GP-modelled successor features (Lemma 1). Additionally, $\gamma \left\| \mathbf{P}_{\mathcal{T}} - \mathbf{P}_j \right\| \times (\left\| \mathbf{Q}_{\mathcal{T}}^* - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_{\mathcal{T}}^j \right\|) / (1-\gamma)$ is a term related to both the dissimilarity of dynamics in environment $j \in \{1, \dots, N\}$ and the target environment \mathcal{T} , as well as the expected future reward in these two environments. We note that Eq. (17) enforces the selection of the best policy among $j \in \{1, \dots, N\}$ policies, hence the derived upper bound is only based on the best selected policy among N source environments.

Theorem 2 is the core of our theoretical analysis that shows by using N different action-value functions that are obtained by their corresponding successor feature functions and reward mapper weights, one can still hold an upper bound on the difference of the optimal policy $Q_{\mathcal{T}}^*(s, a)$, and $\tilde{Q}_{\mathcal{T}}^{\pi}(s, a)$, if $\pi(s)$ is selected by GPI as defined in Eq. (17). As expected, this upper bound is looser when the norm distance between $\tilde{\mathbf{w}}_{\mathcal{T}}$ and \mathbf{w}_j are higher - i.e. the reward functions are significantly different. As explained in Lemma 1, by increasing the amount of noise when modeling the source successor feature functions, $\sigma_{m,d}(\mathbf{x})$ will be higher and accordingly the upper bound will become looser since $\varepsilon(m)$ increases. This is reasonable as increasing $\sigma_{\mathbf{S}}$ indicates that source and target environment are significantly dissimilar. Note that if the environments are exactly similar, $\gamma \left\| \mathbf{P}_{\mathcal{T}} - \mathbf{P}_j \right\| \times (\left\| \mathbf{Q}_{\mathcal{T}}^* - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_{\mathcal{T}}^j \right\|) / (1-\gamma) = 0$, this leads to a similar upper bound in (Barreto et al. (2018)). Further analysis on the obtained upper bound can be found in supplementary materials.

4. Experiments

We evaluate the performance of our method on 3 benchmarks: (1) A toy navigation problem, (2) Classic CartPole control, and (3) The environment introduced in Barreto et al. (Barreto et al. (2020)) and (Barreto et al. (2019)). We compare our approach, Successor Features for Dissimilar Environments (SFDE) with two related studies: Fast Successor Features (FSF) (Barreto et al. (2020)) and Linear Projection of Successor Features (LPSF) (Zhang et al. (2017)). Additional experiments are available in the supplementary materials.

All the algorithms, including SFDE, are used in two phases: (1) The first phase is adaptation where we fine-tune the source successor feature functions using the first 1000 target observations. This step is method specific. (2) A testing phase in which we only use the learnt policy and collect reward without updating the models. For the adaptation phase, ϵ -greedy based exploration is used, afterwards we set the exploration rate $\varepsilon_e = 0$ in the testing phase to demonstrate the effects of transfer from previous environments. The

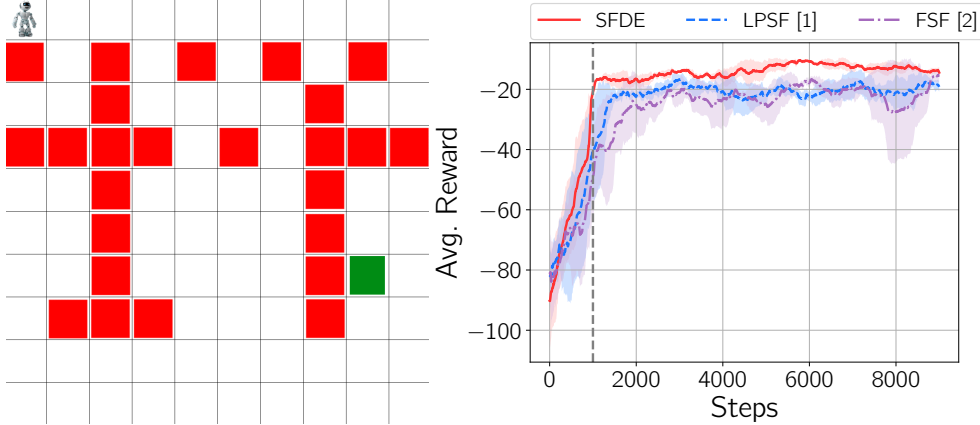


Figure 2. (left) Sample of proposed maze environment with red squares as obstacles with -50 reward and green square as goal with $+100$ reward. (right) Obtained results of SFDE and two baselines. The results are averaged over 10 runs. The dashed vertical line demarcates the adaptation phase.

baseline FSF is not equipped to handle environments with dissimilar dynamics. We adapt this approach by fine-tuning the successor feature functions of the source environments using the first stored 1000 observations of the target environment in the adaptation phase. A batch size of 64 is used at every step $64 \leq t \leq 1000$ to feed the new target observations to the previously learned source successor features. Once the testing phase starts ($t > 1000$), we stop fine-tuning of the successor models for the rest of the experiment. For the LPSF baseline, we follow the idea of (Zhang et al. (2017)) by defining a linear relation between the source and target successor feature - that is there exists a mapping $\beta = \beta_1 \dots \beta_N$ such that the following loss function is minimised: $\mathcal{L}_\beta(\theta_\beta) = \sum_{i=1 \dots N} \left\| \tilde{\psi}^T(\mathbf{x}) - \beta_i \tilde{\psi}^{\mathbf{S}_i}(\mathbf{x}) \right\|, \forall \mathbf{x} \in \mathcal{D}^T$. Intuitively, the best linear projection is found for each source successor feature to minimise its distance to the target successor. We follow the same approach explained for FSF to feed the target observations in adaptation phase to N neural networks that each represent the model of i -th source successor features $i = \{1, \dots, N\}$ with \mathcal{L}_β loss function. Each of these N neural networks are MLPs with no hidden layers that is an equivalent of linear regression in which the weights of the neural networks are β_i . Likewise, the best obtained value of β in the adaptation phase is used in the testing phase. We used the same batch size of 64 to minimise \mathcal{L}_β loss function for LPSF. The linearly projected successor features then used in the GPI framework to obtain the optimal policy. In our approach, we construct the GP based on the combination of source and target observations. To this end, 500 randomly sampled observations from source $\mathcal{D}^{\mathbf{S}}, \forall \mathbf{S} \in \{\mathbf{S}_1, \dots, \mathbf{S}_N\}$ and all the observations in \mathcal{D}^T at every step of the adaptation phase is used. However, once the testing phase is initiated, the new target observations are disregarded and previously seen observations from target are reused. Further details of implementations are

available in the supplementary material.

4.1. Toy Navigation Problem

The proposed maze problem consists of a 10×10 grid and the agent needs to find the goal in the maze. The agent can pass through the obstacles but it receives a reward of -50 . It also receives -1 reward for each step and $+100$ reward for reaching the goal. The action set is defined as $\mathcal{A} = \{\text{left, right, up, down}\}$. Figure 2 (left) shows an example of this environment with red squares indicating obstacles, and the green square as the goal.

A set of 12 policies $\Pi^{\mathbf{S}} = \{\pi^{\mathbf{S}_1}, \dots, \pi^{\mathbf{S}_{12}}\}$ is learnt using generic Q-Learning on randomly generated maze source environments $\{\mathcal{M}(\mathcal{S}, \mathcal{A}, p_1, R_1), \dots, \mathcal{M}(\mathcal{S}, \mathcal{A}, p_{12}, R_{12})\}$ as $\mathbf{S}_1, \dots, \mathbf{S}_{12}$ in which the location of 25 obstacles and goal are changed. After obtaining the corresponding policies, successor feature functions of these source environments are learnt following Algorithm 1. Given $\tilde{\psi}^{\pi_1}, \dots, \tilde{\psi}^{\pi_{12}}$, we generate a random target environment and use the first 1000 observations of adaptation phase as explained. Following Algorithm 2, $\mathcal{GP}^{1 \dots 12}$ is constructed as each step by fitting both source and target observations. We set $\sigma_{\mathbf{S}}^2 = 0.1$ and $\sigma^2 = 0.01$ as the modelling noise and the measurement noise, respectively. Having modelled the successor features, we then perform GPI by using the predicted mean as shown in Eq. (10) and (9). To calculate \tilde{w}_T , the reward mapper of target environment is calculated by minimising loss function introduced in Eq. (6). Figure 2 (right) shows the performance of our approach and other baselines. As expected, our method incorporating the dissimilarity of environments performs better than the other two related approaches. Figure 2 shows that LPSF adjusts the source successor features but it seems to be slower than SFDE in updating the successor feature functions as a linear

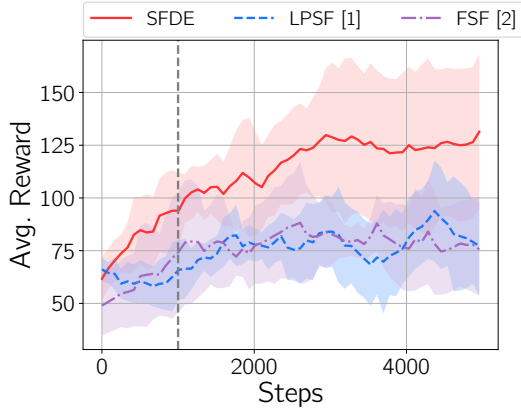


Figure 3. Experiments with CartPole-v0 with $Pole_Length = 0.5m$ in source environment and $Pole_Length = 3m$ in target environment. The results are averaged over 10 runs. The dashed vertical line demarcates the adaptation phase.

projection may not always be found. It can be seen that FSF can update the values of successor features to some extent, however, it seems to be less effective than the other 2 approaches. This can be the result of fine-tuning the successor feature models with significantly different observations that can be an issue in such scenarios. Note that we have not used any visual information in this experiment and the location of agent is translated to (x, y) in the 10×10 grid.

4.2. CartPole-v0

In the CartPole problem, we define a source environment \mathcal{S}_1 with a learnt policy $\pi^{\mathcal{S}_1}$ and a target environment \mathcal{T} . We incorporate the dissimilarities of dynamics in source and target by changing the length of the pole from $Pole_Length = 0.5m$ in source, to $Pole_Length = 3.0m$ in the target environment. This change impacts the transition probabilities of the target environment, hence, it can be considered as a change in dynamics. Similar to the previous experiment, we use the first 1000 observations in all three methods in a same manner. At each step, we fit \mathcal{GP}^1 by combination of source and target observations. We set $\sigma_s^2 = 0.1$, $\sigma^2 = 0.01$, and the maximum number of steps in the CartPole is set to 200 in each episode. To translate the image data into states that can be used in our framework, we used the flattened output of the last convolution layer as the state of the CartPole environment. The detailed structure of the proposed network is available in supplementary materials. Figure 3 shows the results of this experiment and it can be seen that our proposed method outperforms both FSF and LPSF by using the GP-based modelled successor features.

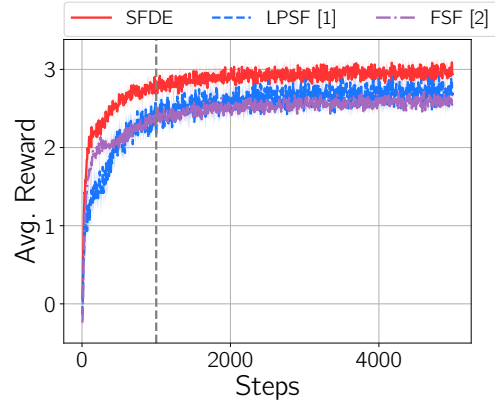


Figure 4. Results on the environment introduced by FSF. The results are averaged over 5 runs. The dashed vertical line demarcates the adaptation phase.

4.3. FSF Environment

We use the environment introduced in Barreto et al. (Barreto et al. (2020)) as our final experiment. The proposed environment is a 10×10 grid cells with $\mathcal{A} = \{\text{left, right, up, down}\}$. There are 10 objects spread across the grid at all time that the agent can pick up. Once an object is picked up, another random object will appear in the grid randomly. Each object belongs to one of two types (red or blue). At each time step, the agent receives an image showing its position, the position of objects, and type of each object (Barreto et al. (2020)). Then the agent selects the proper action to move in a direction. The object is assumed picked up, if the agent occupies the cell in which the object presents. In that case, it gets a reward based on the type of the object and a new object will appear randomly in the grid.

Following the setting of FSF experiments, we trained the agent in 2 different source environments with different reward functions and dynamics of environment. We incorporated the dissimilarity of environments by adding a 5% random transition noise to the target environment and creating a single random terminal state with a negative reward of -1 . Figure 4 shows the obtained results on this problem. Similar to other experiments, our approach seems to outperform the other two baselines. Interestingly, FSF outperforms LPSF in adaptation phase, but both methods approximately converge to the same value of the average reward.

5. Acknowledgements

This research was partially funded by the Australian Government through the Australian Research Council (ARC). Prof Venkatesh is the recipient of an ARC Australian Laureate Fellowship (FL170100006).

6. Conclusion

In this paper we proposed a novel transfer learning approach based on successor features in RL. Our approach is for the scenarios wherein the source and the target environments have dissimilar reward functions as well as dissimilar environment dynamics. We propose the use of Gaussian Processes to model the source successor features functions as noisy measurements of the target successor functions. We provide a theoretical analysis on the convergence of our method by proving an upper bound on the error of the optimal policy. We evaluate our method on 3 benchmark problems and showed that our method outperform existing methods.

References

- Ammar, H. B. and Taylor, M. E. Reinforcement learning transfer via common subspaces. In *International Workshop on Adaptive and Learning Agents*, pp. 21–36. Springer, 2011.
- Barreto, A., Dabney, W., Munos, R., Hunt, J. J., Schaul, T., van Hasselt, H. P., and Silver, D. Successor features for transfer in reinforcement learning. In *Advances in neural information processing systems*, pp. 4055–4065, 2017.
- Barreto, A., Borsa, D., Quan, J., Schaul, T., Silver, D., Hessel, M., Mankowitz, D., Zidek, A., and Munos, R. Transfer in deep reinforcement learning using successor features and generalised policy improvement. In *International Conference on Machine Learning*, pp. 501–510. PMLR, 2018.
- Barreto, A., Borsa, D., Hou, S., Comanici, G., Aygün, E., Hamel, P., Toyama, D., Mourad, S., Silver, D., Precup, D., et al. The option keyboard: Combining skills in reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 13052–13062, 2019.
- Barreto, A., Hou, S., Borsa, D., Silver, D., and Precup, D. Fast reinforcement learning with generalized policy updates. *Proceedings of the National Academy of Sciences*, 117(48):30079–30087, 2020.
- Dayan, P. Improving generalization for temporal difference learning: The successor representation. *Neural Computation*, 5(4):613–624, 1993.
- Gupta, A., Devin, C., Liu, Y., Abbeel, P., and Levine, S. Learning invariant feature spaces to transfer skills with reinforcement learning. *arXiv preprint arXiv:1703.02949*, 2017.
- Joy, T. T., Rana, S., Gupta, S., and Venkatesh, S. A flexible transfer learning framework for bayesian optimization with convergence guarantee. *Expert Systems with Applications*, 115:656–672, 2019.
- Konidaris, G. and Barto, A. Autonomous shaping: Knowledge transfer in reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pp. 489–496, 2006.
- Lederer, A., Umlauf, J., and Hirche, S. Uniform error bounds for gaussian process regression with application to safe control. *arXiv preprint arXiv:1906.01376*, 2019.
- Liang, X., Wang, T., Yang, L., and Xing, E. Cirl: Controllable imitative reinforcement learning for vision-based self-driving. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 584–599, 2018.
- Ma, C., Ashley, D. R., Wen, J., and Bengio, Y. Universal successor features for transfer reinforcement learning. *arXiv preprint arXiv:2001.04025*, 2020.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *nature*, 518(7540): 529–533, 2015.
- Rasmussen, C. E. Gaussian processes in machine learning. In *Summer School on Machine Learning*, pp. 63–71. Springer, 2003.
- Rusu, A. A., Rabinowitz, N. C., Desjardins, G., Soyer, H., Kirkpatrick, J., Kavukcuoglu, K., Pascanu, R., and Hadsell, R. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- Shilton, A., Gupta, S., Rana, S., and Venkatesh, S. Regret bounds for transfer learning in bayesian optimisation. In *Artificial Intelligence and Statistics*, pp. 307–315. PMLR, 2017.
- Siriwardhana, S., Weerasakera, R., Matthies, D. J., and Nanayakkara, S. Vusfa: Variational universal successor features approximator to improve transfer drl for target driven visual navigation. *arXiv preprint arXiv:1908.06376*, 2019.
- Srinivas, N., Krause, A., Kakade, S. M., and Seeger, M. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- Yin, H. and Pan, S. Knowledge transfer for deep reinforcement learning with hierarchical experience replay. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- Zhang, J., Springenberg, J. T., Boedecker, J., and Burgard, W. Deep reinforcement learning with successor features for navigation across similar environments. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2371–2378. IEEE, 2017.

Zhu, Z., Lin, K., and Zhou, J. Transfer learning in deep reinforcement learning: A survey. *arXiv preprint arXiv:2009.07888*, 2020.

7. Supplementary Materials

We show our theoretical analysis for both cases of finite \mathbb{X} and infinite \mathbb{X} .

7.1. Finite \mathbb{X}

Theorem 1

Let \mathbf{S}_i and \mathbf{S}_j be two different source environments with dissimilar transition dynamics p_i and p_j respectively. Let $\delta_{ij} \triangleq \max_{s,a} |r_i(s,a) - r_j(s,a)|$, where $r_i(\cdot, \cdot)$ and $r_j(\cdot, \cdot)$ are the reward functions of environment \mathbf{S}_i and \mathbf{S}_j respectively. We denote π_i^* and π_j^* as optimal policies in \mathbf{S}_i and \mathbf{S}_j . It can be shown that the difference of their action-value functions is upper bounded as:

$$\begin{aligned} Q_i^{\pi_i^*}(s,a) - Q_i^{\pi_j^*}(s,a) &\leq \frac{2\delta_{ij}}{1-\gamma} \\ &+ \frac{\gamma \left\| \mathbf{P}_i(s,a) - \mathbf{P}_j(s,a) \right\|}{(1-\gamma)} \\ &\times \frac{\left(\left\| \mathbf{Q}_i^i - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_i^i \right\| \right)}{(1-\gamma)}, \quad (15) \end{aligned}$$

where $Q_i^{\pi_k^*}$ shows the action-value function in environment \mathbf{S}_i by following an optimal policy that is learned in the environment $\mathbf{S}_k \in \{\mathbf{S}_1, \dots, \mathbf{S}_N\}$. We also define $\mathbf{P}_i(s,a) = [p_i(s'|s,a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{P}_j(s,a) = [p_j(s'|s,a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_i^i = [\max_{b \in \mathcal{A}} Q_i^{\pi_i^*}(s',b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_j^j = [\max_{b \in \mathcal{A}} Q_j^{\pi_j^*}(s',b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_i^j = [\max_{b \in \mathcal{A}} Q_i^{\pi_j^*}(s',b), \dots]_{\forall s' \in \mathcal{S}}$, γ as the discount factor, and $\|\cdot\|$ to be 2-norm (Euclidean norm).

Proof:

We start by following the steps from (Barreto et al. (2018)). The left side of the inequality (15) can be rewritten as:

$$\begin{aligned} Q_i^{\pi_i^*}(s,a) - Q_i^{\pi_j^*}(s,a) &= Q_i^{\pi_i^*}(s,a) - Q_j^{\pi_j^*}(s,a) \\ &+ Q_j^{\pi_j^*}(s,a) - Q_i^{\pi_j^*}(s,a) \\ &\leq \underbrace{Q_i^{\pi_i^*}(s,a) - Q_j^{\pi_j^*}(s,a)}_{(I)} \\ &+ \underbrace{Q_j^{\pi_j^*}(s,a) - Q_i^{\pi_j^*}(s,a)}_{(II)}. \end{aligned}$$

For (I), it can be shown that:

$$\begin{aligned} \left| Q_i^{\pi_i^*}(s,a) - Q_j^{\pi_j^*}(s,a) \right| &\leq \left| r_i(s,a) + \gamma \sum_{s'} p_i(s'|s,a) \right. \\ \max_{s'} Q_i^{\pi_i^*}(s',b) - r_j(s,a) - \gamma \sum_{s'} p_j(s'|s,a) \max_b Q_j^{\pi_j^*}(s',b) \left. \right| &\leq \\ \left| r_i(s,a) - r_j(s,a) \right| + \left| \gamma \sum_{s'} p_i(s'|s,a) \max_b Q_i^{\pi_i^*}(s',b) - \right. \\ &\left. \gamma \sum_{s'} p_j(s'|s,a) \max_b Q_j^{\pi_j^*}(s',b) \right| \leq \\ \left| r_i(s,a) - r_j(s,a) \right| + \gamma \left(\left\| (\mathbf{P}_i - \mathbf{P}_j) \cdot (\mathbf{Q}_i^i - \mathbf{Q}_j^j) \right\| \right) &\leq \\ \left| r_i(s,a) - r_j(s,a) \right| + \gamma \left\| \mathbf{P}_i - \mathbf{P}_j \right\| \times \left\| \mathbf{Q}_i^i - \mathbf{Q}_j^j \right\| = \\ \delta_{ij} + \gamma \left\| \mathbf{P}_i - \mathbf{P}_j \right\| \times \left\| \mathbf{Q}_i^i - \mathbf{Q}_j^j \right\| &\leq \\ \frac{\delta_{ij}}{1-\gamma} + \frac{\gamma}{1-\gamma} \left\| \mathbf{P}_i - \mathbf{P}_j \right\| \times \left\| \mathbf{Q}_i^i - \mathbf{Q}_j^j \right\|, \forall s' \in \mathcal{S} \end{aligned}$$

For (II), it can be shown that:

$$\begin{aligned} \left| Q_j^{\pi_j^*}(s,a) - Q_i^{\pi_j^*}(s,a) \right| &\leq \left| r_j(s,a) + \gamma \sum_{s'} p_j(s'|s,a) \right. \\ \max_b Q_j^{\pi_j^*}(s',b) - r_i(s,a) - \gamma \sum_{s'} p_i(s'|s,a) \max_b Q_i^{\pi_j^*}(s',b) \left. \right| &\leq \\ \left| r_j(s,a) - r_i(s,a) \right| + \left| \gamma \sum_{s'} p_j(s'|s,a) \max_b Q_j^{\pi_j^*}(s',b) - \right. \\ &\left. \gamma \sum_{s'} p_i(s'|s,a) \max_b Q_i^{\pi_j^*}(s',b) \right| \leq \\ \left| r_j(s,a) - r_i(s,a) \right| + \gamma \left(\left\| (\mathbf{P}_j - \mathbf{P}_i) \cdot (\mathbf{Q}_j^j - \mathbf{Q}_i^i) \right\| \right) &\leq \\ \left| r_j(s,a) - r_i(s,a) \right| + \gamma \left\| \mathbf{P}_j - \mathbf{P}_i \right\| \times \left\| \mathbf{Q}_j^j - \mathbf{Q}_i^i \right\| = \\ \delta_{ij} + \gamma \left\| \mathbf{P}_j - \mathbf{P}_i \right\| \times \left\| \mathbf{Q}_j^j - \mathbf{Q}_i^i \right\| &\leq \\ \frac{\delta_{ij}}{1-\gamma} + \frac{\gamma}{1-\gamma} \left\| \mathbf{P}_j - \mathbf{P}_i \right\| \times \left\| \mathbf{Q}_j^j - \mathbf{Q}_i^i \right\| = \\ \frac{\delta_{ij}}{1-\gamma} + \frac{\gamma}{1-\gamma} \left\| \mathbf{P}_i - \mathbf{P}_j \right\| \times \left\| \mathbf{Q}_j^j - \mathbf{Q}_i^i \right\|, \forall s' \in \mathcal{S} \end{aligned}$$

Considering (I) and (II), it leads to the upper bound:

$$\begin{aligned} Q_i^{\pi_i^*}(s,a) - Q_i^{\pi_j^*}(s,a) &\leq \frac{2\delta_{ij}}{1-\gamma} + \frac{\gamma \left\| \mathbf{P}_i - \mathbf{P}_j \right\|}{(1-\gamma)} \\ &\times \frac{\left(\left\| \mathbf{Q}_i^i - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_i^i \right\| \right)}{(1-\gamma)}. \end{aligned}$$

♠

Lemma 1

Let π_1^*, \dots, π_N^* be N optimal policies for $\mathbf{S}_1, \dots, \mathbf{S}_N$ respectively and $\tilde{Q}_{\mathcal{T}}^{\pi_j^*} = (\tilde{\psi}^{\pi_j^*})^\top \tilde{\mathbf{w}}_{\mathcal{T}}$ denote the action-value function of an optimal policy learned in \mathbf{S}_j and executed in target environment \mathcal{T} . Let $\tilde{\psi}^{\pi_j^*}$ denote the estimated successor feature function from the combined source and target observations from \mathbf{S}_j and \mathcal{T} as defined in Eq. (9)(in paper), and $\tilde{\mathbf{w}}_{\mathcal{T}}$ is the estimated reward mapper for environment \mathcal{T} by using loss function in Eq. (6)(in paper). It can be shown that the difference of the true action-value function and the estimated one through successor feature functions and reward mapper, is bounded as:

$$\Pr\left(\left|Q_{\mathcal{T}}^{\pi_j^*}(s, a) - \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(s, a)\right| \leq \varepsilon(m) \forall s, a\right) \geq 1 - \delta,$$

where $\varepsilon(m) = \sqrt{2\log(|\mathbb{X}|u_m/\delta)}\sigma_{m,d}(\mathbf{x})$, $\mathbf{x} \in \mathbb{X}$ $\delta \in (0, 1)$, $u_m = \frac{\pi^2 m^2}{6}$, m being the number of observations in environment \mathcal{T} , and $\mathbf{x} = (s, a)$. $\sigma_{m,d}(\mathbf{x})$ is the square root of posterior variance as defined in Eq. (10)(in paper).

Proof:

For proving this Lemma, we first follow the properties of Normal distribution. Let us assume that $l \sim \mathcal{N}(0, 1)$, and $\psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) \sim \mathcal{N}(\mu_{m,d}(\mathbf{x}), \sigma_{m,d}^2(\mathbf{x}))$, $\mathbf{x} \in \mathbb{X}$, $\mathbf{x} = (s, a)$, and $d = \{1, \dots, D\}$ as defined in Eq. (9) and Eq. (10)(in paper). m target observations are assumed to be available. We follow Srinivas et al. (2009), based on the properties of Normal distribution, $\Pr(l > c)$, $c > 0$ is calculated as:

$$\Pr(l > c) = e^{-c^2/2}(2\pi)^{-1/2} \int e^{-(l-c)^2 - c(l-c)} dl \leq e^{-c^2/2} \Pr(l > 0) = \frac{1}{2}(e^{-c^2/2}).$$

As $c > 0$, we know $e^{-c(l-c)} \leq 1$ for $l \geq c$. Accordingly, $\Pr\left[|\psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) - \tilde{\psi}_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x})| > \beta_m^{1/2}\sigma_{m,d}(\mathbf{x})\right] \leq e^{-\beta_m/2}$, $\beta_m = 2\log(|\mathbb{X}|u_m/\delta)$. Now assuming $l = (\psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) - \tilde{\psi}_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}))/\sigma_{m,d}(\mathbf{x})$ and $c = \beta_m^{1/2}$, the error of modelling successor feature function can be written as:

$$\Pr\left(\left|\psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) - \tilde{\psi}_{\mathcal{T},d}^{\pi_j^*}(s, a)\right| \leq \beta_m^{1/2}\sigma_{m,d}(\mathbf{x})\right) \geq 1 - |\mathbb{X}|e^{-\beta_m/2}. \quad \forall \mathbf{x} \in \mathbb{X}. \quad (16)$$

If $|\mathbb{X}|e^{-\beta_m/2} = \frac{\delta}{u_m}$, the inequality 16 holds for $u_m = \pi^2 m^2/6$. We follow the assumption in (Barreto et al. (2017;

2019; 2020)), $\exists \tilde{\mathbf{w}}_{\mathcal{T}}$, s.t. $\tilde{Q}_{\mathcal{T}}^{\pi_j^*} = (\tilde{\psi}^{\pi_j^*})^\top \tilde{\mathbf{w}}_{\mathcal{T}}$ given all dimensions of successor feature function, hence Lemma 1 holds for $\forall \mathbf{x} \in \mathbb{X}$:

$$\Pr\left(\left|Q_{\mathcal{T}}^{\pi_j^*}(\mathbf{x}) - \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(\mathbf{x})\right| \leq \varepsilon(m) \forall \mathbf{x} \in \mathbb{X}\right) \geq 1 - \delta. \quad \spadesuit$$

We note $\varepsilon(m)$ decreases $\varepsilon(m) \in (\mathcal{O}(\log(m)^{-1}))$ as $\sigma_{m,d}(\mathbf{x}) \in \mathcal{O}(\log(m)^{-2})$ (Lederer et al. (2019)) and $\sqrt{2\log(|\mathbb{X}|u_m/\delta)} \in \mathcal{O}(\log(m))$. This guarantees that the error of modelling convergence to zero as $m \rightarrow \infty$. Before starting the proof of Theorem 2, we present Remark 1 based on the concept of GPI (Barreto et al. (2018)) as follows:

Remark 1 Let π_1, \dots, π_N be N decision policies and correspondingly $\tilde{Q}^{\pi_1}, \tilde{Q}^{\pi_2}, \dots, \tilde{Q}^{\pi_N}$ are the respective estimated action-value functions (Lemma 1) such that:

$$\left|Q^{\pi_i}(\mathbf{x}) - \tilde{Q}^{\pi_i}(\mathbf{x})\right| \leq \varepsilon(m) \quad \forall \mathbf{x} = (s, a),$$

where m is the number of target observations. Defining: $\pi(s) \in \operatorname{argmax}_a \max_i \tilde{Q}^{\pi_i}(s, a)$, Then:

$$Q^\pi(\mathbf{x}) \geq \max_i Q^{\pi_i}(\mathbf{x}) - \frac{2}{1-\gamma}\varepsilon(m) \quad \forall \mathbf{x} = (s, a).$$

Proof:

We start the proof by extending $P^\pi \max_i \tilde{Q}^{\pi_i}(s, a)$, where P^π is the Bellman operator (Barreto et al. (2018)):

$$\begin{aligned} P^\pi \max_i \tilde{Q}^{\pi_i}(s, a) &= r(s, a) + \gamma \sum_{s'} p_i(s'|s, a) \max_i \tilde{Q}^{\pi_i}(s', \pi(s')) \\ &\geq r(s, a) + \gamma \sum_{s'} p_i(s'|s, a) \max_b Q^{\pi_i}(s', b) - \gamma\varepsilon(m) \\ &\geq r(s, a) + \gamma \sum_{s'} p_i(s'|s, a) \max_i Q^{\pi_i}(s', \pi_i(s')) - \gamma\varepsilon(m) \\ &\geq r(s, a) + \gamma \sum_{s'} p_i(s'|s, a) Q^{\pi_i}(s', \pi_i(s')) - \gamma\varepsilon(m) \\ &= Q^{\pi_i}(s, a) - \gamma\varepsilon(m) \end{aligned}$$

as $P^\pi \max_i \tilde{Q}^{\pi_i}(s, a) \geq Q^{\pi_i}(s, a) - \gamma\varepsilon(m)$, $\forall i = \{1, \dots, N\}$:

$$\begin{aligned} P^\pi \max_i \tilde{Q}^{\pi_i}(s, a) &\geq \max_i Q^{\pi_i}(s, a) - \gamma\varepsilon(m) \\ &\geq \tilde{Q}^{\pi_i}(s, a) - \gamma\varepsilon(m) - \varepsilon(m) \\ &\geq \tilde{Q}^{\pi_i}(s, a) - \frac{1+\gamma}{1-\gamma}\varepsilon(m) \\ &\geq Q^{\pi_i}(s, a) - \varepsilon(m) - \frac{1+\gamma}{1-\gamma}\varepsilon(m) \end{aligned}$$

We also know that $Q^\pi(\mathbf{x}) = \lim_{k \rightarrow \infty} (P^\pi)^k \max_i \tilde{Q}^{\pi_i}(\mathbf{x})$.

Then, it follows:

$$Q^\pi(\mathbf{x}) \geq \max_i Q^{\pi_i}(\mathbf{x}) - \frac{2}{1-\gamma} \varepsilon(m) \quad \forall \mathbf{x} = (s, a), \mathbf{x} \in \mathbb{X}.$$

♠

Theorem 2

Let $\mathbf{S}_{i=1 \dots N}$ be N different source environments with dissimilar transition functions $p_{i=1 \dots N}$. Let us denote the optimal policy π that is defined based on the GPI as:

$$\pi(s) \in \operatorname{argmax}_{a \in \mathcal{A}} \max_{j \in \{1 \dots N\}} \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(s, a), \quad (17)$$

where $\tilde{Q}_{\mathcal{T}}^{\pi_j^*} = (\tilde{\psi}^{\pi_j^*})^\top \tilde{\mathbf{w}}_{\mathcal{T}}$ being the action-value function of an optimal policy learned in \mathbf{S}_j and executed in target environment \mathcal{T} , $\tilde{\psi}^{\pi_j^*}$ is the estimated successor feature from the combined source and target observations from \mathbf{S}_j and \mathcal{T} as defined in Eq. (9)(in paper), and $\tilde{\mathbf{w}}_{\mathcal{T}}$ is the estimated reward mapper for target environment from Eq. (6)(in paper). Considering Lemma 1 and Eq. (13)(in paper), the difference of optimal action-value function in the target environment and our GPI-derived action value function is upper bounded as:

$$\begin{aligned} Q_{\mathcal{T}}^*(s, a) - \tilde{Q}_{\mathcal{T}}^{\pi \in \pi_j^*}(s, a) &\leq \frac{2\phi_{\max}}{1-\gamma} \left\| \tilde{\mathbf{w}}_{\mathcal{T}} - \mathbf{w}_j \right\| \\ &+ \frac{\gamma \left\| \mathbf{P}_{\mathcal{T}}(s, a) - \mathbf{P}_j(s, a) \right\|}{(1-\gamma)} \\ &\times \frac{\left(\left\| \mathbf{Q}_{\mathcal{T}}^* - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_{\mathcal{T}}^j \right\| \right)}{(1-\gamma)} \\ &+ \frac{2\varepsilon(m)}{(1-\gamma)}. \end{aligned} \quad (18)$$

where $\phi_{\max} = \max_{s,a} \|\phi(s, a)\|$. We also define $\mathbf{P}_{\mathcal{T}} = [p_{\mathcal{T}}(s'|s, a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{P}_j = [p_j(s'|s, a), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_{\mathcal{T}}^* = [\max_{b \in \mathcal{A}} Q_{\mathcal{T}}^*(s', b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_j^j = [\max_{b \in \mathcal{A}} \tilde{Q}_j^{\pi_j^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, $\mathbf{Q}_{\mathcal{T}}^j = [\max_{b \in \mathcal{A}} \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(s', b), \dots]_{\forall s' \in \mathcal{S}}$, and γ as the discount factor.

Proof:

$Q_{\mathcal{T}}^*(s, a) - \tilde{Q}_{\mathcal{T}}^{\pi \in \pi_j^*}(s, a)$ is defined as the difference of the optimal action-value function, and the action-value function derived from GPI. It can be shown that:

$$\begin{aligned} Q_{\mathcal{T}}^*(s, a) - \tilde{Q}_{\mathcal{T}}^{\pi}(s, a) &\leq Q_{\mathcal{T}}^*(s, a) - Q_{\mathcal{T}}^{\pi_j^*}(s, a) + \frac{2}{1-\gamma} \varepsilon(m) \\ &\leq \frac{2\delta_{ij}}{1-\gamma} + \frac{\gamma \left\| \mathbf{P}_i - \mathbf{P}_j \right\| \times \left\| \mathbf{Q}_{\mathcal{T}}^* - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_{\mathcal{T}}^j \right\|}{(1-\gamma)} \\ &\quad + \frac{2}{1-\gamma} \varepsilon(m) \quad // \text{Theorem 1} \\ &\leq \frac{2}{1-\gamma} \max_{s,a} \|\phi(s, a)\| \|\mathbf{w}_{\mathcal{T}} - \mathbf{w}_j\| \\ &+ \frac{\gamma \left\| \mathbf{P}_i - \mathbf{P}_j \right\| \times \left\| \mathbf{Q}_{\mathcal{T}}^* - \mathbf{Q}_j^j \right\| + \left\| \mathbf{Q}_j^j - \mathbf{Q}_{\mathcal{T}}^j \right\|}{(1-\gamma)} + \frac{2}{1-\gamma} \varepsilon(m) \end{aligned}$$

♠

As explained in Lemma 1, $\varepsilon(m) \rightarrow 0$ with $\varepsilon(m) \in (\mathcal{O}(\log(m))^{-1})$, as $m \rightarrow \infty$ - that is, the number of target observations tend to infinity. Note that the remaining terms of the upper bound depends on the amount of dissimilarity of source and target environments as explained in Section 3.1 of paper.

7.2. Infinite \mathbb{X}

We now continue our analysis on the cases that the action-state space (\mathbb{X}) is infinite - i.e. there may be infinite observations coming from the target environment. In that case, Lemma 1 will not hold and further steps need to be taken.

Let us assume $\mathbb{X}_m \subset \mathbb{X}$ represents a subset of infinite \mathbb{X} at time step m , where m target observations are seen. Clearly, Lemma 1 will hold with this assumption if $\beta_m = 2\log(|\mathbb{X}_m|u_m/\delta)$. The main question in here is if we can extend this to the whole search space \mathbb{X} .

Following Boole's inequality - known as union bound, it can be shown that for some constants $a, b, L > 0$ (Srinivas et al. (2009)):

$$\Pr\left(\forall i = \{1, 2\}, \forall \mathbf{x} \in \mathbb{X}, \left| \frac{\partial \psi_{\mathcal{T},d}^{\pi_j^*}}{\partial \mathbf{x}_i} \right| < L\right) \geq 1 - 2ab \frac{L^2}{b^2},$$

that implies:

$$\left(\forall \mathbf{x} \in \mathbb{X}, \left| \psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) - \psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}') \right| \right) \leq L \|\mathbf{x} - \mathbf{x}'\|, \quad \mathbf{x}, \mathbf{x}' \in \mathbb{X}. \quad (19)$$

Eq. (19) enables us to perform a discretisation on the search space \mathbb{X}_m with size of τ_m^2 so that:

$$|\mathbf{x} - [\mathbf{x}]_m| \leq 2r/\tau_m,$$

where $[\mathbf{x}]_m$ denotes the closest point from \mathbb{X}_m to $\mathbf{x} \in \mathbb{X}$ and τ_m implies the number uniformly spaced points on both coordinates of \mathbb{X}_m that is a discretisation factor. We now proceed to Lemma 2 as an extension of successor feature function modelling error in infinite \mathbb{X} space.

Lemma 2

Let $\mathbf{x} = (s, a) \in \mathbb{X}$, and \mathbb{X} is infinite state-action space. π_1^*, \dots, π_N^* is N optimal policies for $\mathbf{S}_1, \dots, \mathbf{S}_N$ respectively and $\tilde{Q}_{\mathcal{T}}^{\pi_j^*} = (\tilde{\psi}^{\pi_j^*})^\top \tilde{\mathbf{w}}_{\mathcal{T}}$ denote the action-value function of an optimal policy learned in \mathbf{S}_j and executed in target environment \mathcal{T} . Let $\tilde{\psi}^{\pi_j^*}$ denote the estimated successor feature function from the combined source and target observations from \mathbf{S}_j and \mathcal{T} as defined in Eq. (9)(in paper), and $\tilde{\mathbf{w}}_{\mathcal{T}}$ is the estimated reward mapper for environment \mathcal{T} by using loss function in Eq. (6)(in paper). It can be shown that the difference of the true action-value function and the estimated one through successor feature functions and reward mapper, is bounded as:

$$\Pr\left(\left|Q_{\mathcal{T}}^{\pi_j^*}(s, a) - \tilde{Q}_{\mathcal{T}}^{\pi_j^*}(s, a)\right| \leq \varepsilon(m) \forall s, a\right) \geq 1 - \delta,$$

where

$$\varepsilon(m) = \sqrt{2\log(2u_m/\delta) + 8\log(2mbr\sqrt{\log(4a/\delta)})}\sigma_{m,d}([\mathbf{x}]_m) + \frac{1}{m^2}, \mathbf{x} \in \mathbb{X} \delta \in (0, 1), u_m = \frac{\pi^2 m^2}{6}, m > 1 \text{ being the number of observations in environment } \mathcal{T}, [\mathbf{x}]_m \text{ is the closest points in } \mathbb{X}_m \text{ to } \mathbf{x} \in \mathbb{X}. \sigma_{m,d}(\cdot) \text{ is the square root of posterior variance as defined in Eq. (10)(in paper). } a, b > 0 \text{ are constants.}$$

Proof:

As explained in Section 7.2, we know:

$$\Pr\left(\forall \mathbf{x} \in \mathbb{X}, |\psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) - \psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}')| < b\sqrt{\log(4a/\delta)}\right) \geq 1 - \delta$$

Accordingly, by replacing \mathbf{x}' :

$$\forall \mathbf{x} \in \mathbb{X}_m, |\psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) - \psi_{\mathcal{T},d}^{\pi_j^*}([\mathbf{x}]_m)| \leq 2rb\sqrt{\log(4a/\delta)}/\tau_m.$$

By selecting the discretisation factor as $\tau_m = 4m^2br\sqrt{\log(4a/\delta)}$:

$$\forall \mathbf{x} \in \mathbb{X}_m, |\psi_{\mathcal{T},d}^{\pi_j^*}(\mathbf{x}) - \psi_{\mathcal{T},d}^{\pi_j^*}([\mathbf{x}]_m)| \leq \frac{1}{m^2}.$$

This implies $|\mathbb{X}_m| = (4m^2br\sqrt{\log(4a/\delta)})^2$. By replacing $|\mathbb{X}_m|$ in β_m defined in Section 7.2, the proof is completed. \spadesuit

Hence, if \mathbb{X} is infinite set, Theorem 2 holds with $\varepsilon(m) = \sqrt{2\log(2u_m/\delta) + 8\log(2mbr\sqrt{\log(4a/\delta)})}\sigma_{m,d}([\mathbf{x}]_m) + \frac{1}{m^2}$.

7.3. Experimental Details

Maze: For the task of navigation, we designed a maze environment with following properties: (1) We set the ε -greedy exploration rate to $\varepsilon_e = 0.5$ for the adaptation phase with decay rate of 0.9999, this value is set to zero in the testing phase. (2) Discount factor value is set to $\gamma = 0.9$. (3) $\alpha = 0.05$ is the learning rate. Agent is allowed to reach to the goal in maximum of 100 steps, otherwise it terminates.

For 12 source environments, 25 obstacles are randomly generated, the agent always start from top left, and the goal is also randomly placed in these environments. We used generic Q-learning with replay buffer size 10^4 and Adam optimizer with batch size 64 to find the optimal policies in all these 12 environments. Algorithm 1 in the paper is then used to estimate the successor feature functions in the environment. Our proposed feature function for this problem is a MLP with 4 hidden layers with a linear activation function in the last hidden layer to represent the reward mapper of the task. The remaining hidden layers have ReLU activation function. The output of this network is the predicted value of the reward for an state and action. Note that this network minimises the loss function introduced in Eq. (6)(in paper).

CartPole: As mentioned, to translate the image data into states, we used a CNN with: (1) First hidden layer with 64 filters of 5×5 with stride 3 with a ReLU activation, second hidden layer with 64 filters of 4×4 with stride 2 and a ReLU activation, third hidden layer with 64 filters of 3×3 with stride 1 and a ReLU activation. The final hidden layer is ‘‘features’’ we used in an image that is fully connected

(20) Maximum number of steps for CartPole problem is set to 200. We set the ε -greedy exploration rate to $\varepsilon_e = 0.5$ for the adaptation phase with decay rate of 0.9999, this value is set to zero in testing phase. The source learned policy is with pole’s length of $0.5m$ and accordingly, using Algorithm 1 in the paper the corresponding successor features are extracted. We used generic Q-learning with replay buffer size 10^5 and Adam optimizer with batch size 64 to find this optimal policy. The target environment is then modified to incorporate the change of environment. Figure 5 demonstrates the change of dynamics. We used the same structure of feature function in Maze problem for this experiment.



Figure 5. Illustration of the change in the dynamics for the CartPole problem. (Left) Pole’s length is $0.5m$ in the source environment and (Right) Pole’s length changed to $3m$ in the target environment.

FSF: This environment (Barreto et al. (2020)) is a 10×10 grid with 10 objects and an agent occupying one cell at each time step. There are 2 types of objects each with a reward associated with it. We randomly initialise those 10 objects by sampling both their type and position from a uniform distributions over the corresponding sets. Likewise, the initial position of the agent is a uniform sample of all possible positions in the grid. The reward function is defined by the object type that has been picked up by the agent. e.g. Picking up red object is $+1$ reward and picking up blue is -1 . Agent picks up an object if it occupies that particular cell in which the object exists. If agent picks up an object, another one will be generated randomly (in terms of location and type) in the grid. At each step the agent receives an observation representing the configuration of the environment (Barreto et al. (2020)). These are $11 \times 11 \times (\mathbb{D} + 1)$ tensors that can be seen as 11×11 images with $(\mathbb{D} + 1)$ channels that are used to identify objects and walls (Barreto et al. (2020)). The observations are shifted so that the the agent is always at the top-left cell of the grid. Figure 6 shows an example of this environment. The 2 source policies used in this experiment have $w_1 = [1, 0]$, $w_2 = [1, 0]$, respectively. Intuitively, the reward mappers indicate picking up an object of particular type and ignoring the other type. However, in the target environment, the change of reward function is to pick up the first object type and “avoid” the second one with negative reward - i.e. $w_{\mathcal{T}} = [1, -1]$. For the dissimilarity of dynamics, we added 5% noise to the transitions of the agent and also randomly placed a terminal state in the target environment with -1 reward.

7.4. Additional Experiments

In this section, we compare our results in navigation problem with generic Q-Learning. We relaxed the assumption adaptation and testing phase and Q-Learning is allowed to use all the observations from the target environment. Figure 7 shows Q-Learning also converges to the same amount of avg. reward, however, since no transfer is involved, the it is significantly slower than other baselines. For Q-Learning, we set $\varepsilon_e = 0.9$ with a decay rate of 0.9999 in 10^4 iterations.

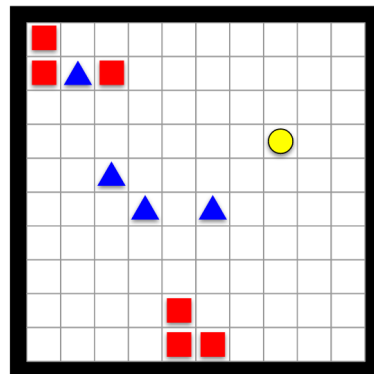


Figure 6. An example of the environment described in FSF (Barreto et al. (2020)).

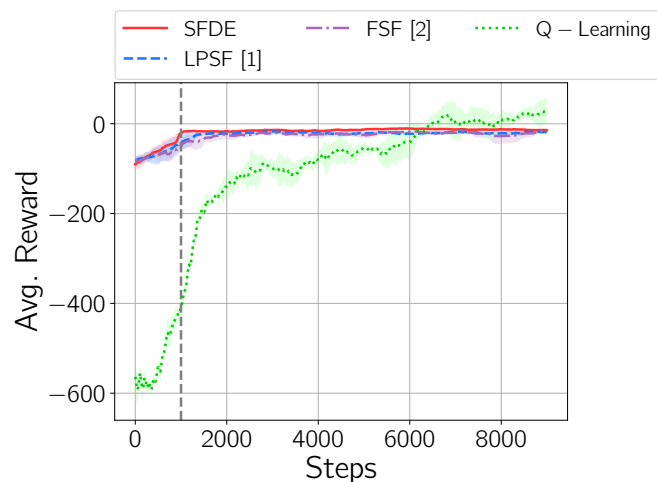


Figure 7. Performance of Q-Learning in navigation problem. [1] and [2] references are shown in the paper.