

Reinforcement Learning Transfer using a Sparse Coded Inter-Task Mapping

Haitham Bou Ammar¹, Matthew E. Taylor², Karl Tuyls¹, and Gerhard Weiss¹

¹ Department of Knowledge Engineering, Maastricht University
The Netherlands

² Department of Computer Science, Lafayette College
USA

Abstract. Reinforcement learning agents can successfully learn in a variety of difficult tasks. A fundamental problem is that they learn slowly in complex environments, inspiring the development of speedup methods such as transfer learning. Transfer improves learning by reusing learned behaviors in similar tasks, usually via an inter-task mapping, which defines how a pair of tasks are related. This paper proposes a novel transfer learning technique to autonomously construct an inter-task mapping by using a novel combinations of sparse coding, sparse projection learning, and sparse pseudo-inputs gaussian processes. Experiments show successful transfer of information between two very different domains: the mountain car and the pole swing-up task. This paper empirically shows that the learned inter-task mapping can be successfully used to (1) improve the performance of a learned policy on a fixed number of samples, (2) reduce the learning times needed by the algorithms to converge to a policy on a fixed number of samples, and (3) converge faster to a near-optimal policy given a large amount of samples.

1 Introduction

Reinforcement Learning (RL) is a popular framework that allows agents to solve sequential-action selection tasks with minimal feedback. Unfortunately, RL agents may learn slowly in large or complex environments due to the computational effort needed to attain an acceptable performing policy. Transfer Learning [16] (TL) is one technique used to cope with this difficulty by providing a good starting prior for the RL agent attained in a related source task.

The source task can thus differ from the target task in many ways. If the tasks have different representations of state or action spaces, some type of mapping between the task is required. This inter-task mapping matches each state/action pair of the source task to its corresponding state/action pair in the target facilitating transfer. While there have been a number of successes in using such a mapping, the approaches are typically hand-coded and may require substantial human knowledge [18, 16]. Our contributions in this paper are twofold. First, we propose a novel scheme to automatically learn an inter-task mapping between two tasks. Second, we introduce the new *Transfer Least Squares Policy Iteration* (TrLSPI) algorithm for transfer between tasks of continuous state spaces and discrete action spaces.

To the best of our knowledge, this paper shows the first successful attempts to automatically transfer between two very different RL benchmarks. Namely, we conduct experiments to automatically transfer from the Mountain Car to the Pole Swing-up

problem. Our results show (1) improved performance on a fixed number of samples, (2) a reduction in the convergence times to attain a policy on a fixed number of samples, and (3) a reduction in the time needed to attain a near-optimal policy on a large amount of samples.

The rest of the paper proceeds as follows. Related work is discussed next in Section 2. Background information is presented in Section 3. Section 4 describes how an inter-task mapping can be learned between two tasks by leveraging sparse coding, sparse projection learning and sparse pseudo-input gaussian processes. In Section 5, we introduce our novel TrLSPI algorithm showing how the learned mapping can be used to transfer information between a source task and target task. Experiments of transfer between two very different tasks is presented in Section 6. Section 7 presents a discussion on the scope and applicability of our framework. Section 8 concludes with a discussion of future work.

2 Related Work

In the past few years there has been a significant amount of work done in transfer learning for RL tasks. This section outlines the most related work and contrasts it with the work in this paper.

The majority of current transfer learning work in RL assumes that either 1) the two agents are very similar and no mapping is needed, or 2) the inter-task mapping is provided by a human. For instance, [18] transfers advice and [16] transfers Q-values — both methods assume that a mapping between the state and action variables in the two tasks has been provided. Another approach is to frame different tasks as having a shared *agent space* [4], so that no explicit mapping is needed. However, this requires that the agent acting in both tasks share the same actions and a human to map new sensors back into the agent space. The primary contrast between these methods and the current work is that we are interested in *learning* a mapping between states and actions in pairs of tasks, rather than assuming that it is provided or unnecessary.

Our previous work [1] required the presence of hand-coded shared features to automatically learn the inter-task mapping. This work extends the previous approach to overcome the need for a *predefined common subspace* to determine the inter-task mapping.

There has been recent work that approaches fully autonomous transfer. For example, semantic knowledge about state features between two tasks may be used [8, 5], background knowledge about the range or type of state variables can be used [13, 17], or transition models for each possible mapping could be generated and tested [14]. Transfer learning has also been successful across different domains, e.g., using a simple discrete, deterministic task to improve learning on a complex, continuous, noisy task [15]. However, there are currently no general methods to learn an inter-task mapping without requiring (1) background knowledge, which is not typically present in RL settings, or (2) an expensive analysis of an exponential number of inter-task mappings. This paper overcomes these problems by automatically discovering high level features and using them to conduct transfer within reasonable time requirements.

Unlike all other existing methods (to the best of our knowledge) and complementary to our previous work [1, 14, 15], we assume differences among all the variables of Markov Decision Processes describing the source and target tasks, and focus on learning

an *inter-state mapping*, rather than a state-variable mapping. Additionally, our framework can use state-dependent action mappings, allowing flexibility that other existing algorithms do not.

3 Background

This section provides the reader with a short overview of Reinforcement Learning, Gaussian Processes, Sparse Coding, Transfer Learning and other learning methods used in this paper.

3.1 Reinforcement Learning (RL)

In a RL problem, an agent must decide how to sequentially select actions to maximize its expected long term reward [12, 2]. Such problems are typically formalized as Markov decision processes (MDPs). An MDP is defined by $\langle S, A, P, R, \gamma \rangle$, where S is the (potentially infinite) set of states, A is the set of all possible actions that the agent may execute, $P : S \times A \rightarrow S$ is a state transition probability function describing the transition dynamics, $R : S \rightarrow \mathbb{R}$ is the reward function measuring the performance of the agent, and $\gamma \in [0, 1)$ is the discount factor. A policy $\pi : S \rightarrow A$ is defined as a mapping from a state to an action. The goal of an RL agent is to improve its policy, potentially reaching the optimal policy π^* :

$$Q^*(s, a) = \max_{\pi} E\left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \mid s = s_0, \pi\right] \quad (1)$$

In tasks with continuous state and/or action spaces, the Q functions and/or policies cannot be represented in a table format, typically requiring sampling or function approximation techniques. This paper uses one such technique, *Least Squares Policy Iteration* (LSPI), which will be explained next.

3.2 Transfer Learning in RL Tasks

Typically, when using Transfer Learning (TL) in RL tasks, there are a source and target task [16]. When the source and the target tasks are related, transferring a learned source behavior should improve learning speed in the target task by providing an informative prior. The prior will restrict the exploration in the target task by biasing the agent so that it chooses actions that are better than random exploration, reducing the target task learning times and improving the overall performance. In our formulation, each of these tasks is defined as an MDP which is a tuple of $\langle S^{(i)}, A^{(i)}, P^{(i)}(s, a), R^{(i)}, \gamma^{(i)} \rangle$ for $i \in \{1, 2\}$ where $S^{(i)}$, $A^{(i)}$, $P^{(i)}(s, a)$, $R^{(i)}$ and $\gamma^{(i)}$ represent the state spaces, action spaces, state transition probabilities, reward functions and discount factors for each of the source ($i = 1$) and target ($i = 2$) tasks.

The source and the target task may differ in their state spaces and/or action spaces (as well as other components of the MDP). If transfer is to be useful when such differences exist, an inter-task mapping relating these state-action spaces differences [16] can be used. Our main focus in this paper is the automatic discovery of the inter-task mapping that enables transfer. The upcoming sections will further clarify the need for such a mapping as well as describe our novel framework in attaining it.

3.3 Least Squares Policy Iteration

LSPI [6] is an approximate RL algorithm that is considered an actor/critic method. LSPI is composed of two parts. The first is an evaluation step, *Least Squares Temporal Difference Q-learning* (LSTDQ) and the second is a policy improvement step. LSTDQ is an evaluation step: the algorithm will update the weights representing the policy so that the new parameters minimize certain error criteria. For example, the LSTDQ could be set to minimize the Bellman residual error of the projected Bellman equations. Once this step has finished, LSPI uses the attained weights to improve the policy by taking greedy actions in the approximated Q -function.

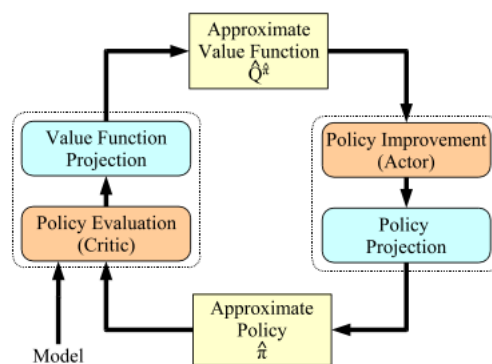


Fig. 1. Least Squares Policy Iteration schematic [6]

Figure 1 highlights the actor and critic organization of LSPI. Since LSPI uses function approximators to represent the Q -functions and/or policies, there also exist two projection phases for both the Q -function and the policy, as can be seen in the schematic. A thorough mathematical treatment may be found elsewhere [6].

3.4 Gaussian Processes

Gaussian Processes (GPs) form a research field by themselves. It is beyond the scope of this paper to fully detail the mathematical framework. This section briefly explains GPs and refers the reader elsewhere [10] for a more in-depth treatment.

Parametric models have traditionally been used to solve regression problems because of their interpretability. These parametric models may face problems when trained on complex data sets that require highly expressive models. Since the inter-task mapping is considered to be an expressively complex relation, GPs were preferred over simple parametric forms.

GPs are a form of supervised learning used to discover a relation between a given set of input vectors, \mathbf{x} , and output pairs, \mathbf{y} . As opposed to normal regression techniques that perform inference in the weight space, GPs perform inference directly in the functional space, making learning simpler. Since a function can be represented as an infinite dimensional vector and GPs are distributions over functions, a GP is an extension of a

multidimensional Gaussian distribution into infinite dimensions. Following existing notation [10], if a function is sampled according to a GP we write:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')), \quad (2)$$

where $m(\mathbf{x})$ and $k(\mathbf{x}, \mathbf{x}')$, represent the mean and covariance function that fully specify a GP.

Learning in a GP setting involves maximizing the marginal likelihood:

$$\log p(\mathbf{y}|\mathbf{x}) = -\frac{1}{2}\mathbf{y}^T \mathbf{K}^{-1} \mathbf{y} - \frac{1}{2} \log |\mathbf{K}| - \frac{n}{2} \log 2\pi. \quad (3)$$

Maximizing Equation 3 may be computationally complex as we must invert the covariance matrix \mathbf{K} , which is of order of $O(N^3)$, where N is the number of input points. Therefore, we use a fast learning technique, *Sparse Pseudo-input Gaussian Processes* (SPGP), as proposed elsewhere [11].

3.5 Sparse Pseudo Input Gaussian Processes

SPGPs aim to reduce the complexity of learning and prediction in GPs by parametrizing the regression model with $M \ll N$ pseudo-input points, while still preserving the full Bayesian framework. The covariance of the GP model is parametrized by the location of the $M \ll N$ pseudo-inputs. Mathematically, the marginal likelihood to be maximized is defined as:

$$p(\mathbf{y}|\mathbf{X}, \bar{\mathbf{X}}) = \mathcal{N}(\mathbf{y}, \mathbf{K}_{NM} \mathbf{K}_M^{-1} \mathbf{K}_{MN} + \mathbf{\Lambda} + \sigma^2 \mathbf{I}), \quad (4)$$

$\bar{\mathbf{X}}$ is the matrix formed by the pseudo-inputs with $\bar{\mathbf{X}} = \{\bar{\mathbf{x}}\}_{m=1}^M$. \mathbf{K}_{NM} the covariance matrix formed by the pseudo-inputs and the real inputs as $\mathbf{K}_{MN} = k(\bar{\mathbf{x}}_m, \mathbf{x}_n)$ with $k(\cdot, \cdot)$ being the covariance kernel. \mathbf{K}_M^{-1} is the inverse of the covariance matrix formed among the pseudo inputs with $\mathbf{K}_M = k(\bar{\mathbf{x}}_m, \bar{\mathbf{x}}_m)$. \mathbf{K}_{MN} is the covariance matrix formed by the pseudo-inputs and the real inputs. $\mathbf{\Lambda}$ is a diagonal matrix having the diagonal entries $\lambda_n = k_{nn} - \mathbf{k}_n^T \mathbf{K}_M^{-1} \mathbf{k}_n$. The noise variance and the identity matrix are represented by σ and \mathbf{I} respectively.

Existing results [11] show a complexity reduction in the training cost (i.e., the cost of finding the parameters of the covariances) to $O(M^2N)$ and in the prediction cost (i.e., prediction on a new set of inputs) to $O(M^2)$. The results further demonstrate that the SPGPs framework can fully match normal GPs with small M (i.e., few pseudo-inputs), successfully producing very sparse solutions. A full mathematical treatment may be found elsewhere [11].

3.6 Sparse Coding

Sparse Coding (SC) [7] is an unsupervised learning technique used to find a high level representation for a given set of unlabeled input data. It does this by discovering a succinct over-complete basis for the provided data set. Given m k -dimensional vectors, ζ , SC aims to find a set of n basis vectors, \mathbf{b} , and activations, a , with $n > k$ such that $\zeta \approx \sum_{j=1}^n a_j^{(i)} \mathbf{b}_j$, where i and j represent the number of input data patterns and number of bases, respectively. SC begins by assuming a Gaussian and a sparse prior on the reconstruction error ($\zeta^{(i)} - \sum_{j=1}^n a_j^{(i)} \mathbf{b}_j$) and on the activations, solving the following an optimization problem:

$$\begin{aligned}
\min_{\{\mathbf{b}_j\}, \{a_j^{(i)}\}} & \sum_{i=1}^m \frac{1}{2\sigma^2} \|\zeta^{(i)} - \sum_{j=1}^n \mathbf{b}_j a_j^{(i)}\|_2^2 \\
& + \beta \sum_{i=1}^m \sum_{j=1}^n \|a_j^{(i)}\|_1 \\
s.t. & \|\mathbf{b}_j\|_2^2 \leq c, \forall j = \{1, 2, \dots, n\}
\end{aligned} \tag{5}$$

The problem presented in Equation 5 is considered to be a “hard” optimization problem as it is not jointly convex (i.e., in the activations and bases). However, fast and efficient optimization algorithms exist [7] and were used in our work, as described in Section 4.1.

4 Learning the Inter-Task Mapping

In order to automatically construct an inter-task mapping, χ , relating two different MDPs, we propose a novel framework utilizing sparse coding, a L_1 projection scheme, and sparse pseudo-input Gaussian processes. Each of these methods is necessary to solve a problem that is inherent to TL in RL tasks. We approach the problem of learning the inter-task mapping, χ , as a supervised learning problem. As χ is a mapping relating state-action triplets from the source with the target task, related triplets should be provided as training data points. Unfortunately, this is itself a hard problem — it is not trivial for the user to describe what state triplets in the source task correspond to what in the target task. We therefore approach this problem by automatically transforming the problem spaces (i.e., the state-action spaces of the two tasks) into a higher representational space through SC, projecting the target task data onto those attained bases and then utilizing a Euclidean distance measure to gauge similarity (Section 4.2). At this stage, the data set is ready and is to be provided to the regressor to construct the inter-task mapping. Many regression techniques could be applied to the approach but we chose to use a non-parametric approximation scheme because of its generalization advantages.

The following sections further clarify each of the above steps and explain the technicalities involved.

4.1 The need for an Inter-Task Mapping

For transfer between different but related MDPs to be possible, an inter-task mapping is essential [16]. The source and target MDPs may differ in the state and/or action spaces. If the transfer scheme is to be successful, there must be a mapping function that can overcome the differences of the MDPs. Traditionally, such a mapping was thought to be a one-to-one mapping between the state/action variables representing the tasks [16]. Differing from the traditional way, we think of such a mapping to be a function that relates state-action successor state triplets from the source with the target task. In other words, our inter-task mapping is more than just a one-to-one mapping between the state and/or action spaces of the MDPs. It also includes other terms that are automatically discovered by our global approximators, which ultimately enhances the transfer approach.

Mathematically, $\chi : S_s \times A_s \times S_s \rightarrow S_t \times A_t \times S_t$, where S and A represent the state space and the action space of the source and the target task, respectively.

The inter-task learning framework can conceptually be split into three essential parts. The first is the dimensional unification of both the source and target task state-action spaces of the MDPs. The second is the automatic discovery of a high dimensional informative space of the source task. This is achieved through SC, as described in Section 4.2, ensuring that transfer is conducted in a high representational space of the source task. In order to use a similarity measure among different patterns, the data should be present in the same space. That is why the target task samples still need to be projected to the attained high representational space of the source. This is done using sparse projection learning, described in Section 4.3. The third and final step is to approximate the inter-task mapping via a non-parametric regression technique, explained in Section 4.4.

4.2 Sparse Coding Transfer for RL

As described in Section 3.6, SC is an efficient way to discover higher level information in a given unlabeled data set. We use SC to solve two inherent problems in transfer learning for RL tasks. The first is to unify the dimensions of the state action spaces of the two different MDPs. The second is to discover a higher level representation for the attained bases and activations of the source task state-action spaces. This step guarantees that our scheme works with the “best” available representation/information space of the source task.

Unifying the Source and Target Dimensions Our problem commences by first unifying the dimensions of the state action spaces of the two MDPs, an essential step for discovering the inter-task mapping. After this step has finished, any existing TL in RL technique may be used. However, this paper goes further and proposes a new transfer framework based on the attained bases and activations, described in Section 5.

This “dimensional unification” process is described in Algorithm 1. In short, Algorithm 1 sparse codes random samples from the source task constrained by learning the same number of bases (d_t) as those of the target task.

The algorithms proposed elsewhere [7] solve Equation 6 in line 3 of Algorithm 1. After this stage is done, new activations and bases describing the samples are attained.³ Note that, these newly attained samples—described as a linear combinations of the bases and activations (i.e. \mathbf{Ab})—do not yet relate anything to the target task ones. The target task samples still need to be projected towards these bases. This is done as described in Section 4.3.

After Algorithm 1 is finished, new features in the source task state action spaces are discovered. This is reasonable as TL typically transfers between a low dimensional source task to a high dimensional target task. Here, SC is determining new bases that are of a higher number than the original state action dimensions in the source task. If successful, new patterns and representations are discovered in the source task state-action spaces. These new features describe new representations not anticipated by the

³ Please note that while writing Algorithm 1 it was assumed that the dimensions of the source task d_s are lower than those of the target task d_t . But it is worth noting that it works as well for the other cases.

Algorithm 1 Sparse Coding Transfer Reinforcement Learning

Require: Source MDP samples $\{\langle s_s, a_s, s'_s \rangle\}_{i=1}^m$, Target MDP samples $\{\langle s_t, a_t, s'_t \rangle\}_{j=1}^f$

- 1: Calculate d_s and d_t which are the dimensions of each of the state action spaces of the MDPs
- 2: Sparse code the source by solving:
- 3:

$$\begin{aligned} \min_{\{\mathbf{b}_j\}, \{a_j^{(i)}\}} \sum_{i=1}^m \frac{1}{2\sigma^2} \|\langle s_s, a_s, s'_s \rangle^{(i)} - \sum_{j=1}^{d_t} b_j a_j^{(i)}\|_2^2 & \quad (6) \\ & + \beta \sum_{i=1}^m \sum_{j=1}^{d_t} \|a_j^{(i)}\|_1 \\ \text{s.t. } \|b_j\|_2^2 \leq c, \forall j = \{1, 2, \dots, d_t\} \end{aligned}$$

- 4: Solve the problem of Equation 6 using the algorithm proposed in [7]
 - 5: Return the Activation matrix ($\mathbf{A} \in \mathbb{R}^{m \times d_t}$) and the Bases ($\mathbf{b} \in \mathbb{R}^{d_t \times 1}$)
-

original dimensions. Therefore, this new information can be used to help and guide the transfer learning scheme.

High Information Representation After dimensional unification, as described in the previous section, SC is again used to discover a succinct higher informational/representational bases of the activations than the unified dimensional spaces. This insures that our transfer approach operates in the “richest” space described through the samples. This is done in a similar framework to that in Section 4.2 and is described in Algorithm 2.

Algorithm 2 Succinct High Information Representation of MDPs

Require: Activations acquired through Algorithm 1, Number of new high dimensional bases d_n

- 1: Represent the activations in the d_n bases by solving the following problem using the algorithm in [7]:
- 2:

$$\begin{aligned} \min_{\{\mathbf{z}_j\}, \{c_j^{(i)}\}} \sum_{i=1}^m \frac{1}{2\sigma^2} \|\langle \mathbf{a}_{1:d_t} \rangle^{(i)} - \sum_{j=1}^{d_n} \mathbf{z}_j c_j^{(i)}\|_2^2 & \quad (7) \\ & + \beta \sum_{i=1}^m \sum_{j=1}^{d_n} \|c_j^{(i)}\|_1 \\ \text{s.t. } \|\mathbf{z}_j\|_2^2 \leq o, \forall j = \{1, 2, \dots, d_n\} \end{aligned}$$

- 3: **return** Return new activations $\mathbf{C} \in \mathbb{R}^{m \times d_n}$ and bases $\mathbf{z} \in \mathbb{R}^{d_n \times 1}$
-

The idea presented by Algorithm 2 is to sparse code the activations, representing the original samples of the MDPs, to a higher representational space, d_n .⁴ This stage should

⁴ In our experiments we have set d_n to be 100, a relatively high number.

guarantee that we project the samples of the source task MDP into a high informational space where a similarity measure can be used to find a relation between the source and target task triplets. Noting that there are no restrictions on the number of bases to be determined: unneeded bases have an activation of zero once the SC problem has been solved.

At this stage, the source state action spaces are described in a rich informational space determined by the newly discovered bases and activations. The next step is to project the target task samples to that space described by \mathbf{Z} so that triplet can be ordered and the inter-task mapping approximated.

4.3 L_1 Sparse Projection Learning

Once the above stages have finished, the source samples are described via the activations generated in Algorithm 2. However, target task samples still have no relationship to the learned activations. In other words, the bases and activations that have been attained successfully describe high informational patterns and representations in the source task state-action spaces but do not represent the target state-action spaces. Since we are seeking a similarity correspondence between the source and target task triplets, the target task samples should be represented in the high informational space of the source task.

Therefore, the next step is to learn a sparse projection to project the target task samples onto the \mathbf{Z} basis representing the source task MDP. In other words, the goal now is to learn a sparse projection that is capable of representing the random target task samples as a combination of some activations, automatically learned, and the \mathbf{Z} bases generated by Algorithm 2. The overall scheme is described in Algorithm 3, where the activations are learned by solving the L_1 regularized least squares optimization problem of Equation 8. This optimization problem guarantees that the attained activations are as sparse as possible and is solved using the interior point method [3].

At this stage all the samples from both the target and source task are projected to the same space described by the sparse coded vectors \mathbf{Z} . The next step will be to order the data points from both the source and the target task so to approximate the inter-task mapping.

Algorithm 3 Reflecting Target Task Samples

Require: Sparse Coded Bases \mathbf{Z} generated by Algorithm 2, Target MDP samples

$$\{(s_t, a_t, s'_t)\}_{i=1}^f$$

1: **for** $i = 1 \rightarrow f$ **do**

2: Represent the target data patterns in the sparse coded bases, \mathbf{Z} , by solving:

3:

$$\hat{\phi}^{(i)}(\langle s_t, a_t, s'_t \rangle) = \arg \min_{\phi^{(i)}} \|\langle s_t, a_t, s'_t \rangle - \sum_{j=1}^{d_n} \phi_j^{(i)} z_j\|_2^2 + \beta \|\phi^{(i)}\|_1 \quad (8)$$

4: **end for**

5: **return** Activations Φ

4.4 Similarity Measure and Inter-Task Mapping Approximation

As mentioned previously, we tackle the problem of learning an inter-task mapping via supervised learning. Since χ maps triplets from the source task to their corresponding triplets in the target task, the problem at this stage is to attain the training patterns to approximate χ .

After reaching the rich space representing the random samples of the 2 MDPs (i.e., \mathbf{Z}), an Euclidean distance measure is used to compare triplets, providing a data set to the regressor (i.e, SPGPs) to approximate the inter-task mapping χ . This similarity measure is used to determine the correspondence of the source and target tasks triplets. Once applied, the similarity measure will seek the triplets of the source task closest to those of the target task and map them together as being inputs and outputs for the regression algorithm, respectively. This is shown in line 2 of Algorithm 4. Since the similarity measure is used in the sparse coded spaces, the distance is calculated using the attained activation (\mathbf{C} and Φ) rather than the samples themselves. Therefore, the scheme has to trace the data back to the original dimensions of the state-action pairs of the MDPs.

There are few restrictions on the function approximation techniques to could be used. We use nonparametric regression and with Sparse Gaussian Processes technique [11]. We prefer Sparse Gaussian Processes rather than normal Gaussian Processes regression technique as the latter may have problems dealing with large data sets. To better clarify, consider the learning phase of a GP that involves minimizing Equation 3. It is clear that the inversion of the covariance matrix, \mathbf{K} , is required on each iteration with complexity $O(n^3)$, where n is the number of samples. Additionally, the minimization algorithm (Conjugate Gradient Descent [9]) may get stuck in a local minimum of Equation 3, a common problem in function approximation schemes and minimization problems.

Algorithm 4 Similarity Measure & Inter-Task mapping approximation

Require: Sparse Coded Basis \mathbf{Z} , Sparse Coded Activations of the source task $\mathbf{C} \in \mathbb{R}^{m \times d_n}$,
Projected Target Task activations $\phi \in \mathbb{R}^{m \times d_n}$

- 1: **for all** ϕ **do**
 - 2: Calculate the closest activation in \mathbf{C} minimizing the Euclidean/similarity distance measure.
 - 3: **end for**
 - 4: Correspond the triplets with the minimum similarity measure as being inputs and outputs to create a data set \mathcal{D}
 - 5: Approximate the Inter-task mapping, χ using SPGPs
 - 6: **return** The approximated Inter-task mapping χ
-

5 Transfer Scheme

Assuming there exists a “good enough” policy, π_s^* for the source task, we propose a novel transfer algorithm for pairs of tasks with continuous state spaces and discrete action spaces, titled Transfer Least Squares Policy Iteration.

This section describes the novel transfer scheme and reflects on the details and technicalities of the approach.

5.1 Transfer Least Squares Policy Iteration

Algorithm 5 can be split into two sections. The first determines χ (of Section 4), using source task samples⁵ (from π_s^*). The second provides those samples for the evaluation phase of the LSPI algorithm (LSTDQ) as a start, to learn on and improve the policy on the target task. The intuition here is that if the tasks are similar and if the inter-task mapping is “good enough”, then those samples will bias the target task controller towards choosing good actions and restricting its area of exploration thus reducing learning times and increasing performance.

Algorithm 5 TrLSPI

Require: Source MDP samples $\{\langle s_s, a_s, s'_s \rangle\}_{i=1}^n$, Target MDP samples $\{\langle s_t, a_t, s'_t \rangle\}$, Number for re-samples n_s , close to optimal policy for the source system π_s^* , State action basis functions for the target task ψ_1, \dots, ψ_k

- 1: *Unify the Dimensions* using Algorithm 1
- 2: Discover *High Informational Representation* using Algorithm 2
- 3: *Sparse Project* the target task samples using Algorithm 3
- 4: Use a *similarity* measure to attain the data set and *approximate* χ using Algorithm 4
- 5: Randomly Sample n_s Source task triplets $\langle s_s, a_s, s'_s \rangle_{i=1}^{n_s}$ greedily in the optimal policy π_s^* , set of state-dependent basis function $\psi_1, \dots, \psi_k : S_t \times A_t \rightarrow \mathbb{R}$
- 6: **for** $i = 1 \rightarrow n_s$ **do**
- 7: Find the corresponding target task triplets as $\langle s_t^{(i)}, a_t^{(i)}, s_t^{(i)'} \rangle = \chi(\langle s_s^{(i)}, a_s^{(i)}, s_s^{(i)'} \rangle)$
- 8: **end for**
- 9: Find the closest triplet in the initial samples to the ones predicted by χ
- 10: Use LSTDQ described by [6] to evaluate those samples
- 11: Learn and Improve Policy till convergence using LSPI [6]
- 12: **return** Learned Policy π_{target}^*

Provided that the tasks are related, Algorithm 5 is capable of attaining a good starting behavior for the target task. The performance of this policy depends on the state space region where those samples were provided. In other words, it is not possible to achieve near-optimal performance with a small number of samples that are in regions far from the goal state.⁶ Therefore, if the agent has to seek a near-optimal policy, then either a new sampling step using the current policy should be added to Algorithm 5, or a large amount of samples should be provided. It is worth noting that it is not necessary for the algorithm to be provided by a model for the system to perform that sampling. A black box generative model taking inputs being states and actions and producing outputs of successor states and rewards is sufficient.

6 Experiments & Results

Two very different tasks were chosen to evaluate the proposed framework, the RL benchmark tasks Mountain Car (MC) and Pole Swing-Up (see Figure 4).

⁵ If using an approximate RL algorithm in the source task, the policy would not be optimal rather a near-optimal one which works as well.

⁶ This is a problem that is inherit to LSPI and not to our TrLSPI algorithm.

The control objective of MC, the source task, is to drive the car up the hill (Figure 2). The difficulty is that gravity is stronger than the car’s motor—even at maximum throttle the car can not directly reach the top of the hill. The solution is to first move away from the target to the opposite side of the hill and then accumulate enough energy to reach the top of the hill. The dynamics of the car are described via two continuous state variables (x, \dot{x}) representing the position and velocity of the center of gravity of the car, respectively. The input action takes on three distinct values: maximum throttle forward (+1), zero throttle (0), and maximum throttle reverse (-1). The car is rewarded by +1 once it reaches the top of the hill, -1 if it hits the wall, and zero elsewhere.

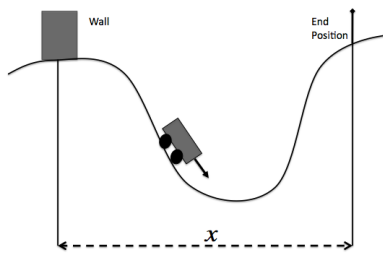


Fig. 2. Mountain Car

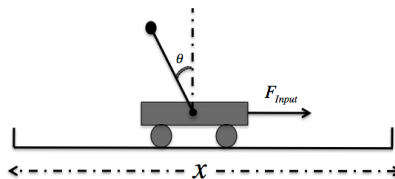


Fig. 3. Pole Swing-Up

Fig. 4. Mountain Car to Pole Swing-up Transfer

The target task is the Pole Swing-up problem described in Figure 3. The control goal of the pole swing-up system is balancing the pole in an upright position (i.e., $\theta = \dot{\theta} = 0$). The allowed actions are (+1) for full throttle right and (-1) for full throttle left. The reward function of the system consisted of two parts: (1) $\cos(\theta)$, which yields its maximum value of +1 at the upright position of the pole, and (2) -1 if the cart hits the boundaries of the track. The angle was restricted to be within $|\theta| < \frac{\pi}{9}$ while the position was restricted to $|x| < 3$.

As clear from the description, the two MDPs representing the tasks are significantly different. The source and target task have different state spaces, action spaces, transition probabilities, and reward functions. No previous work can learn to autonomously transfer between such different tasks.

Our framework requires an optimal policy in the MC source task, π_{MC}^* . SARSA(λ) [12] is used to learn π_{MC}^* . The policy, once learned, was then used to randomly sample different numbers of initial states of task, to be used by χ . We started with 5000 and 2000 randomly sampled states (using a random policy) for the Mountain Car and the Pole Swing-up, respectively. These samples were used by the algorithm described in Section 4 to attain the inter-task mapping χ . After χ has been learned, different amounts of samples were sampled from the source task using the optimal policy π_{MC}^* . Specifically, we have sampled 500, 1000, . . . 20000 states as input to the TrLSPI algorithm to measure performance and convergence times.

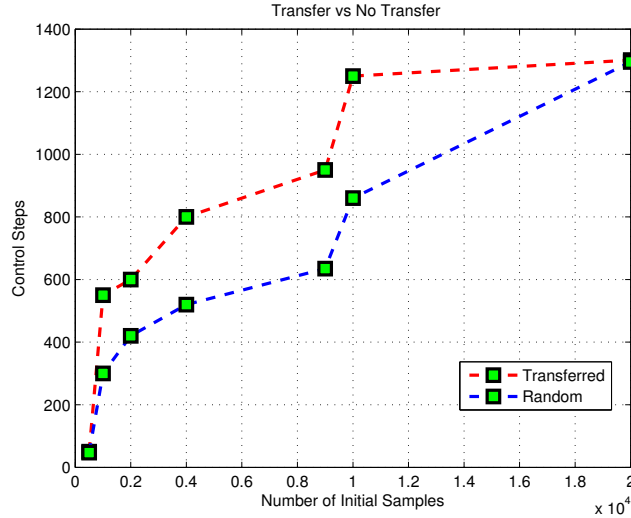


Fig. 5. Transfer Results on the Pole Swing-up task.

Our results show both an increase in the performance on a fixed number of samples and a decrease in the convergence times in both a predefined number of samples and to attain an optimal policy. We measured the performance as the number of successful steps to control the pole in an upright position on a given fixed amount of samples. Another performance measure was the convergence times of the algorithm to a policy on a given fixed number of samples and to learn a near-optimal policy of the target task. The graph of Figure 5 summarizes the results attained applying our framework on different number of transferred samples and compares them with those attained through normal LSPI learning scheme. The graph in Figure 5 clearly shows an increase in the number of control steps in the case of the transferred samples compared to a random sampling scheme. For example it can be seen that at a small number of samples, e.g. 2000, our transfer scheme was able to attain an average of 600 control steps with about 400 for the random case. This performance increases with the number of samples to reach 800 steps at 4000 transferred samples. The random case needed to be provided by 9000 samples to attain such an 800 steps performance. Finally, the transferred algorithm and the random selection scheme seem to converge, on a large amount of samples 20000 to the same number of control steps, about 1300. This leads to the following conclusion:

Conclusion 1: TrLSPI has provided a better distribution of samples compared to random policy in the target task.

The other improvement we report is the decrease in the convergence times, represented by the number of iterations in LSPI, provided a fixed amount of transferred sampled. To better clarify, LSPI was able to converge faster once provided the transferred samples compared to a random sample data set. For example, it took LSPI 7 iterations to converge provided 5000 transferred samples with 12 iterations for the random case. Further the algorithm converged within 14 iteration provided 20000 transferred samples while it took it about 21 for the random case.

Conclusion 2: TrLSPI converged faster provided a fixed amount of samples.

Finally, LSPI was able to converge to an acceptable policy within a 22.5 minutes after being provided a random data set, compared to 17 minutes with the transferred data set⁷. Calculating χ took an additional 3.7 minutes.

Conclusion 3: TrLSPI converged faster to sub-optimal policy compared to a random selection scheme.

7 Discussion & Scope of the Framework

We speculate that the framework we propose is applicable to any model-free TL in RL problem with continuous state spaces and discrete action spaces, covering many real world RL problems. The framework has the advantage of automatically finding the inter-task functional mapping using SC and any “good” regression technique. But there is one potential weakness, as discussed next.

Our framework should work correctly when the two tasks at hand are *semantically* similar, as the rewards of the two systems were not taken into account in the explained scheme. For instance, consider the transfer example between the same robot but with opposite reward tasks. In other words, the robot have the same transitions in the two tasks but have to reach two opposite states.

Our mapping scheme of Section 4, once applied, will produce a one-to-one mapping from the source to the target task. In other words, since the two tasks have the same state and action spaces, the mapping that will be a one-to-one, mapping the same state action successor state triplets of the two tasks together. Therefore, the transition of the robots to the rewardable/ un-rewardable states will map together. Since the optimal policies of the two robots are opposite, it is easy to see that in this case the target task has been provided with a “bad” biased starting policy which will worsen the agents performance rather than enhancing it. Another name for this problem is negative transfer, which is a well known problem in TL for RL tasks. We think that our approach will be able to avoid this scheme once the rewards are added to the similarity measure generating the training set to approximate the inter-task mapping χ . The improvement of this measure to incorporate the rewards is out of the scope of this paper and will be looked at in details in our future work.

8 Conclusions & Future Work

This paper has presented a novel technique for transfer learning in reinforcement learning tasks. Our framework may be applied to pairs of reinforcement learning problems with continuous state spaces and discrete action spaces. The main contributions of this paper are (1) the novel method of automatically attaining the inter-task mapping, χ and (2) the new TrLSPI algorithm for tasks with continuous state spaces and discrete actions. We approached the problem by framing the approximation of the inter-task mapping as a supervised learning problem that was solved using Sparse Pseudo Input Gaussian Processes. Sparse Coding, accompanied with a similarity measure, was used to determine the data set required by the regressor for approximating χ . Our results demonstrate successful transfer between two very different tasks, the mountain car to the pole swing-up task. Success was measured both in an increase in learning performance as well as a reduction in convergence time. We speculate that the process usefully

⁷ Our experiments were performed on a 2.8 Ghz Intel Core i7

restricts exploration in the target task and that the transferred state quality resulting from our scheme.

There are many exciting directions for future work. First, we will compare different distance metrics and demonstrate their effects on the overall performance of the algorithm. Second, the distance measure will be improved by incorporating the rewards in the framework, helping to avoid the problem of negative transfer, as well as reflect upon a criterion for TL in RL.

References

1. H. B. Ammar and M. E. Taylor. Common subspace transfer for reinforcement learning tasks. In *Proceedings of the Adaptive and Learning Agents workshop (at AAMAS-11)*, May 2011.
2. L. Buşoniu, R. Babuška, B. De Schutter, and D. Ernst. *Reinforcement Learning and Dynamic Programming Using Function Approximators*. CRC Press, Boca Raton, Florida, 2010.
3. S. Jean Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky. An interior-point method for large-scale ℓ_1 -regularized logistic regression. *Journal of Machine Learning Research*, 2007, 2007.
4. G. Konidaris. A framework for transfer in reinforcement learning. In *Proceedings of the ICML-06 Workshop on Structural Knowledge Transfer for Machine Learning*, 2006.
5. G. Kuhlmann and P. Stone. Graph-based domain mapping for transfer learning in general games. In *Proceedings of The Eighteenth European Conference on Machine Learning*, September 2007.
6. M. G. Lagoudakis and R. Parr. Least-squares policy iteration. *J. Mach. Learn. Res.*, 4:1107–1149, December 2003.
7. H. Lee, A. Battle, R. Raina, and A. Y. Ng. Efficient sparse coding algorithms. In *In NIPS*, pages 801–808. NIPS, 2007.
8. Y. Liu and P. Stone. Value-function-based transfer for reinforcement learning using structure mapping. In *Proceedings of the Twenty-First National Conference on Artificial Intelligence*, pages 415–20, July 2006.
9. J. Nocedal and S. J. Wright. *Numerical optimization*. Springer, Aug. 1999.
10. C. E. Rasmussen. In *Gaussian processes for machine learning*. MIT Press, 2006.
11. E. Snelson and Z. Ghahramani. Sparse gaussian processes using pseudo-inputs. In *Advances In Neural Information Processing Systems*, pages 1257–1264. MIT press, 2006.
12. R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*, 1998.
13. E. Talvitie and S. Singh. An experts algorithm for transfer learning. In *Proceedings of the Twentieth International Joint Conference on Artificial Intelligence*, 2007.
14. M. E. Taylor, G. Kuhlmann, and P. Stone. Autonomous transfer for reinforcement learning. In *Proceedings of the Seventh International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 283–290, May 2008.
15. M. E. Taylor and P. Stone. Cross-domain transfer for reinforcement learning. In *Proceedings of the Twenty-Fourth International Conference on Machine Learning (ICML)*, June 2007.
16. M. E. Taylor and P. Stone. Transfer learning for reinforcement learning domains: A survey. *J. Mach. Learn. Res.*, 10:1633–1685, December 2009.
17. M. E. Taylor, S. Whiteson, and P. Stone. Transfer via inter-task mappings in policy search reinforcement learning. In *Proceedings of the Sixth International Joint Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, pages 156–163, May 2007.
18. L. Torrey, T. Walker, J. Shavlik, and R. Maclin. Using advice to transfer knowledge acquired in one reinforcement learning task to another. In *In Proceedings of the Sixteenth European Conference on Machine Learning*, pages 412–424, 2005.