
Deep Transfer for Model-Free Reinforcement Learning Using Autonomous Intertask Mappings and Q-Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper...

2 1 Introduction

3 Transfer Learning (TL) within the Reinforcement Learning (RL) domain can be described as lever-
4 aging mastery of one task (a source task) to improve learning speed or asymptotic performance in
5 another task (a target task). The solution to achieving effective transfer depends on the differences
6 between the two task goals, the environments the tasks are posed in, and the agents trying to solve
7 them. The particular problem of deep transfer is the realm of transfer problems where the two tasks
8 may differ in all three of these dimensions.

9 Reinforcement learning in a realistic setting usually does not allow for either *a priori* knowledge of
10 the environment dynamics or *a priori* knowledge of the relationship between observations, actions,
11 and observed rewards; usually agents must rely solely on experience when learning how to master an
12 RL problem. The method of transfer in RL utilizing actual recorded source task experience is known
13 as instance transfer, and is especially difficult in the deep transfer space. Deep instance transfer can
14 be achieved using a linear or non-linear mapping from the source experience space to the target
15 experience space, and for fully autonomous intelligent transfer to occur this mapping must be learned
16 by the target task agent independent of any explicit human mapping. With an inter-task mapping an
17 agent can translate source task experience into pseudo target task experience and use these pseudo
18 experiences as initial samples in a sample-based RL algorithm such as fitted Q-learning.

19 2 Reinforcement Learning

20 3 Transfer Learning for Reinforcement Learning

21 [introduce TL broadly]

22 Instance Transfer

23 Shallow Instance Transfer

24 Deep Instance Transfer

25 4 Restricted Boltzmann Machines

26 4.1 Bipartite RBMs

27 4.2 High-Order RBMs

28 5 Related Work

29 5.1 Early Work

30 *See Taylor review - Early Taylor

31 5.2 Contrasting Methods

32 Gupta Progressive Nets *See newer review

33 5.3 Taylor’s MASTER Algorithm

34 [explain method and why it works] [explain limitations, motivating Ammar’s work]

35 5.4 Ammar’s TrRBM Method

36 The main contribution of this paper is to build off the dissertation work of Ammar (CITE) in which
37 he uses a high-order Restricted Boltzmann Machine to learn intertask mappings for instance transfer.
38 Ammar’s method uses a 3-way RBM with one layer for each task’s concatenated instance tuple
39 vectors (s,a,s?) and one hidden layer to modulate interaction between the two visible layers. The
40 visible layer nodes $\mathbf{v}_1, \mathbf{v}_2$ are modelled as Gaussian random variables, $v_1^{(i)} \sim \mathcal{N}(\mu^{(i)}, \sigma), v_2^{(j)} \sim$
41 $\mathcal{N}(\mu^{(j)}, \sigma)$ and the hidden layer nodes \mathbf{h} take the value of sigmoidal activations. Ammar gives two
42 formulations for the 3-way RBM; the full Transfer Restricted Boltzmann Machine (TrRBM) version
43 with a 3-way weight tensor having elements \mathcal{W}_{ijk} , and a factored ?fTrRBM? (Factored Transfer
44 Restricted Boltzmann Machine) version in which the 3-way weight tensor is factored into the product
45 of 3 layer-specific matrixes, . The factored version is motivated by a need to reduce computational
46 complexity from the full version’s $O(n^3)$ to a more manageable complexity of $O(n^2)$. [FIGURE for
47 fTrRBM] [Talk about why the TrRBM can learn a good mapping] [Talk about learning in the TrRBM
48 model i.e. mean of gaussians] [Motivate the extensions i.e. black-box model is unrealistic, sampling
49 method is unrealistic]

50 6 New Extensions

51 This section is for showing our formalisms for the extensions we are making

52 6.1 TrRBM for a Model-Free Setting

53 6.1.1 Using Source Q-Values Instead of Black-Box Target Task Rewards

54 6.1.2 Transferring Best- and Worst- Policy Instances

55 6.1.3 More Realistic Initial Sampling

56 7 Experimental Setup

57 7.1 Environments

58 2D Mountain Cart

59 3D Mountain Cart

60 2D Cartpole

61 **3D Cartpole**

62 **Acrobot**

63 **2D Maze**

64 **3D Maze**

65 **Breakout**

66 **Pong**

67 **7.2 Untested Design Choices**

68 Briefly explain all parameter/model choices, why we did not tune/experiment with these, and what
69 effect these might have!

70 **8 Results**

71 Display individual experiments tables comparing modifications against baselines and each other?

72 Display individual plots showing the same?

73 **9 Limitations**

74 **10 Conclusions**

75 **References**