Improving Action Selection in MDP's via Knowledge Transfer

Alexander Sherstov and Peter Stone

Department of Computer Sciences The University of Texas at Austin

Overview

- Learning setting
 - Must learn *multiple* tasks in the same domain
 - Actions not uniformly relevant
 - Designed for *large* action sets
- Solution: action transfer
 - Usually beneficial (pastry chef, driver, bidding agent)
 - Formalism + analysis of action transfer
 - Enhancement: *randomized task perturbation*
 - Empirical validation

MDP Formalism



- S is a set of **states**, A is a set of **actions**
- $\mathbf{t}: \mathcal{S} \times \mathcal{A} \to \Pr(\mathcal{S})$ is a transition function
- $r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is a *reward function*
- **policy** is any mapping $\mathcal{S} \to \mathcal{A}$; want policy with maximum expected return in all states

Running Example: Grid World Domain



- states = cells, actions = { $(d, p) : d \in \{\uparrow, \downarrow, \rightarrow, \leftarrow\}, p \in [0.5, 0.9]$ }
- move succeeds w/ prob. proportional to p, random o/w
- reward: -1/2 in quicksand, 1/2 in goal, $-p^2$ o/w
- optimal actions have $p \in [0.5, 0.6]$

Need for Related-Task Formalism



- r, t defined in terms of state space \mathcal{S}
- r, t incomparable across tasks
- But: want to exploit *abstract* transition/reward dynamics

Eliminating S: **Outcomes**

• Original definition:

$$\mathbf{t}: \mathcal{S} \times \mathcal{A} \to \Pr(\mathcal{S})$$
$$r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$$

• Definition using **outcomes** \mathcal{O} :

$$\mathbf{t}: \mathcal{S} \times \mathcal{A} \to \Pr(\mathcal{O})$$

 $r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$

• Example in grid world: $\mathcal{O} = \{\uparrow, \downarrow, \rightarrow, \leftarrow, \mathsf{STAY}\}.$

Eliminating S: **Classes**

• Previous definition:

$$\mathbf{t}: \mathcal{S} \times \mathcal{A} \to \Pr(\mathcal{O})$$
$$r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$$

• Final definition, using outcomes \mathcal{O} and *classes* \mathcal{C} :

$$\mathbf{t}: \mathcal{C} \times \mathcal{A} \to \Pr(\mathcal{O})$$
$$r: \mathcal{C} \times \mathcal{A} \to \mathbb{R}$$

• Example in grid world: $C = \{ EMPTY, GOAL, QUICKSAND \}$.

Complete Formalism



$$\mathbf{t}: \mathcal{C} \times \mathcal{A} \to \Pr(\mathcal{O})$$

 $r: \mathcal{C} \times \mathcal{A} \to \mathbb{R}$

- Similarities across tasks are implicit, unstated.
- Dissimilarity metric $\Delta(U, \tilde{U})$ *re-expresses* these high-level similarities in a *precise, analytically tractable* geometric quantity.
- $\Delta(U, \tilde{U})$ expressed i.n.o. the new formalism and V^*

Bound Based on Task Similarity

Transfer results in a value drop of at most

 $\Delta(U, \tilde{U}) \cdot \sqrt{2\gamma}/(1-\gamma)$

at each state.

Reducing Task Dissimilarity

- Optimization possible if $\mathcal{O}, \mathcal{C}, \kappa$, and η known
- Alternative: *uniform sampling* of value space Complexity:

$$\Omega\left(rac{(v_{\max}-v_{\min})^n}{\epsilon^n}
ight)$$
 draws.

 More informed search: randomized task perturbation (RTP)

RTP Action Transfer at Work



RTP Transfer in Pseudocode (Q**-Learning)**

- $new \rightarrow 1$ Add each $s \in S$ to \mathcal{F} with probability ϕ
 - $\rightarrow 2$ foreach $s \in \mathcal{F}$
 - $\begin{array}{lll} \rightarrow 3 & \text{do } random\text{-}value \leftarrow \operatorname{rand}(v_{\min}, v_{\max}) \\ \rightarrow 4 & Q^+(s, a) \leftarrow random\text{-}value \text{ for all } a \in \mathcal{A} \\ 5 & \text{repeat } s \leftarrow \text{current state}, a \leftarrow \pi(s) \\ 6 & \text{Take action } a, \text{observe reward } r, \text{state } s' \\ 7 & Q(s, a) \xleftarrow{\alpha} r + \gamma \max_{a' \in \mathcal{A}} Q(s', a') \\ \rightarrow 8 & \text{if } s \in S \setminus \mathcal{F} \\ & \text{then } Q^+(s, a) \xleftarrow{\alpha} r + \gamma \max_{a' \in \mathcal{A}} Q^+(s', a') \\ 9 & \text{until converged} \end{array}$

$$\rightarrow 10 \quad \mathcal{A}^* = \bigcup_{s \in \mathcal{S}} \{ \arg \max_{a \in \mathcal{A}} Q(s, a) \}$$
$$\rightarrow 11 \quad \mathcal{A}^+ = \bigcup_{s \in \mathcal{S} \setminus \mathcal{F}} \{ \arg \max_{a \in \mathcal{A}} Q^+(s, a) \}$$

 \rightarrow 12 return $\mathcal{A}^* \cup \mathcal{A}^+$

Relevance-weighted action selection

- RELEVANCE $(a) = |\{s \in \mathcal{S} : \pi^*(s) = a\}|/|\mathcal{S}|.$
- Choice probability on *exploratory* moves proportional to relevance
- Major performance gains

Empirical Results: Grid World



Action Sets

• Full:



• Optimal (value iteration):



Action Sets, cont.

• Transferred:



• RTP-transferred (sample run):



Performance



Related Work

- Transfer in MDP's:
 - hierarchical (Hauskrecht et al., 1998, Dietterich, 2000),
 - first-order (Boutilier et al., 2001),
 - factored (Guestrin et al., 2003)
- Limitation: reliance on description of similarities
- RTP: no guidance, robust to noise, focuses on *actions*

Summary

Contributions:

- Theoretical abstraction for transfer learning
- Formal analysis + transfer quality guarantees
- Empirical validation

Future work:

- Combine with non-transfer approaches to action selection
- Enable fully continuous learning

References

- (Boutilier et al., 2001) Boutilier, C., Reiter, R., and Price, B. (2001). Symbolic dynamic programming for first-order MDPs. In *Proc. 17th International Joint Conference on Artificial Intelligence (IJCAI-01)*, pages 690–697, Seattle, WA.
- (Dietterich, 2000) Dietterich, T. G. (2000). Hierarchical reinforcement learning with the MAXQ value function decomposition. *Journal of Artificial Intelligence Research*, 13:227–303.
- (Guestrin et al., 2003) Guestrin, C., Koller, D., Gearhart, C., and Kanodia, N. (2003). Generalizing plans to new environments in relational MDPs. In *Proc. 18th International Joint Conference on Artificial Intelligence (IJCAI-03)*.
- (Hauskrecht et al., 1998) Hauskrecht, M., Meuleau, N., Kaelbling, L. P., Dean, T., and Boutilier, C. (1998). Hierarchical solution of Markov decision processes using macroactions. In *Proc. Fourteenth Conference on Uncertainty in Artificial Intelligence (UAI-98)*, pages 220–229.

Definition 1. A domain is a quintuple $\langle \mathcal{A}, \mathcal{C}, \mathcal{O}, \mathbf{t}, r \rangle$, where \mathcal{A} is a set of actions; \mathcal{C} is a set of state classes; \mathcal{O} is a set of action outcomes; $\mathbf{t} : \mathcal{C} \times \mathcal{A} \to \Pr(\mathcal{O})$ is a transition function; and $r : \mathcal{C} \times \mathcal{A} \to \mathbb{R}$ is a reward function.

Definition 2. A *task* within the domain $\langle \mathcal{A}, \mathcal{C}, \mathcal{O}, \mathbf{t}, r \rangle$ is a triple $\langle \mathcal{S}, \kappa, \eta \rangle$, where \mathcal{S} is a set of states; $\kappa : \mathcal{S} \to \mathcal{C}$ is a state classification function; and $\eta : \mathcal{S} \times \mathcal{O} \to \mathcal{S}$ is a next-state function.

Definition 3. The outcome value vector of state s in the task $\langle S, \kappa, \eta \rangle$ within the domain $\langle A, C, O, \mathbf{t}, r \rangle$ is the vector $[V^*(s_1) \quad V^*(s_2) \quad \dots \quad V^*(s_{|\mathcal{O}|})]^{\mathrm{T}}$, where $V^* : S \to \mathbb{R}$ is the optimal value function of the task, and each $s_i = \eta(s, o_i)$ is a successor state of s upon outcome $o_i \in O$.

Definition 4. Let $U = \langle U_{c_1}, \ldots, U_{c_{|C|}} \rangle$ and $\tilde{U} = \langle \tilde{U}_{c_1}, \ldots, \tilde{U}_{c_{|C|}} \rangle$ be the OVV sets of the primary and auxiliary tasks, respectively. The **dissimilarity** of the primary and auxiliary tasks, denoted $\Delta(U, \tilde{U})$, is:

 $\Delta(U, \tilde{U}) \stackrel{\text{def}}{=} \max_{c \in \mathcal{C}} \max_{\mathbf{u} \in U_c} \left\{ \min_{\tilde{\mathbf{u}} \in \tilde{U}_c} ||\mathbf{u} - \tilde{\mathbf{u}}||_2 \right\}.$

Detrimental Action Transfer



Auxiliary task



Primary task



Auxiliary Task:

$[3 \textcircled{4} \textcircled{4} 3 3]^{\mathrm{T}}$	$[4 (5) (5) 3 4]^{T}$	$[5 \bigcirc 6 4 5]^{\mathrm{T}}$	[6 ⑦	6	5	$[6]^{T}$
$[35544]^{T}$	$[\ 4\ \textcircled{6}\ \textcircled{6}\ 4\ 5\]^{\mathrm{T}}$	$\begin{bmatrix} 5 (7) & (7) 5 6 \end{bmatrix}^{\mathrm{T}}$	[6 (9)	7	6	$[7]^{\mathrm{T}}$
$[\ 4\ \textcircled{6}\ \textcircled{6}\ 5\ 5\]^{\mathrm{T}}$	$[\ 5\ \widehat{(7)}\ \widehat{(7)}\ 5\ 6\]^{\mathrm{T}}$	$[6 \ 9 \ 9 \ 6 \ 7]^{\mathrm{T}}$	[7(10)	9	7	$[9]^{T}$
$\begin{bmatrix} 5 & 6 & 7 & 6 & 6 \end{bmatrix}^{\mathrm{T}}$	$[\begin{array}{ccc} 6 & 7 \\ \end{array} \begin{array}{c} 9 \\ 6 \\ 7 \end{array} \begin{array}{c} 7 \end{array} \left[\begin{array}{c} 7 \\ \end{array} \right]^{\mathrm{T}}$	$\left[\begin{array}{cc}7&9 ext{(10)}7&9\end{array} ight]^{\mathrm{T}}$	[(10 10 (10	10	$(10)^{T}$

Primary Task:

$[79977]^{T}$	$\begin{bmatrix} 9 (10) & 7 & 7 & 9 \end{bmatrix}^{\mathrm{T}}$	$\begin{bmatrix} 7 \ 9 \ 6 \ 9 \ 7 \end{bmatrix}^{\mathrm{T}}$	$[6 (7) 6 (7) 6]^{T}$
$\left[\begin{array}{cc} 7 & 7 \end{array} \right]^{\mathrm{T}}$	$[\begin{array}{ccc} \textcircled{10} \begin{array}{ccc} \textcircled{10} \begin{array}{ccc} \textcircled{10} \begin{array}{ccc} \textcircled{10} \begin{array}{ccc} \textcircled{10} \end{array} \end{array} \end{matrix}]^{\mathrm{T}}$	$\left[\begin{array}{cc} 7 & 7 & 7 & 10 \\ \end{array} 9 \right]^{\mathrm{T}}$	$[\begin{array}{ccc} 6 & 6 & 7 \\ 9 & 7 \end{bmatrix}^{\mathrm{T}}$
$\left[\begin{array}{ccc} 9 & 6 & \begin{array}{c} 9 & 7 \end{array} \right]^{\mathrm{T}}$	$[\begin{array}{cccccccccccccccccccccccccccccccccccc$	$[\ 9 \ 6 \ 6 \ 9 \ 7]^{\mathrm{T}}$	$[\bigcirc 5 \ 6 \bigcirc 6]^{\mathrm{T}}$
$\left[\begin{array}{ccc} \hline 7 & 6 & \hline 7 & 6 & 6 \end{array}\right]^{\mathrm{T}}$	$[9 7 6 6 7]^{\mathrm{T}}$	$\left[\begin{array}{cccc} \hline 0 & 6 & 5 & \hline 0 & 6 \end{array}\right]^{\mathrm{T}}$	$[\bigcirc 55 \bigcirc 5]^{\mathrm{T}}$