

AVALIAÇÃO DE POLÍTICAS ABSTRATAS NA TRANSFERÊNCIA DE CONHECIMENTO EM NAVEGAÇÃO ROBÓTICA

COMPARISON AND EVALUATION OF ABSTRACT POLICIES FOR TRANSFER LEARNING IN ROBOT NAVIGATION TASKS

Rafael Lemes Beirig, Fernando Andrade Pereira, Marcelo Li Koga, Tiago Matos, Valdinei Freire da Silva, Anna Helena Reali Costa

rafaelbeirigo@usp.br, fernando.pereira2@poli.usp.br, mlk@usp.br, tiago.matos@poli.usp.br, valdinei.freire@usp.br, anna.reali@poli.usp.br

Universidade de São Paulo, Brasil.

Abstract: This paper presents a new approach to the problem of solving a new task by the use of previous knowledge acquired during the process of solving a similar task in the same domain, robot navigation. A new algorithm, Qab-Learning, is proposed to obtain the abstract policy that will guide the agent in the task of reaching a goal location from any other location in the environment, and this policy is compared to the policy derived from another algorithm, ND-TILDE. The policies are applied in a number of different tasks in two environments. The results show that the policies, even after the process of abstraction, present a positive impact on the performance of the agent.

Keywords: Relational representation; Abstract policy; Reinforcement learning; Inductive logic programming; Robotics.

Resumo: Este paper apresenta uma análise comparativa do desempenho de políticas abstratas na reutilização de conhecimento no domínio de navegação robótica. O algoritmo Qab-Learning é proposto para a obtenção dessas políticas, que são então comparadas a políticas abstratas geradas por um algoritmo presente na literatura, o ND-TILDE, na solução de diferentes tarefas de navegação robótica. O presente trabalho demonstra que essas políticas, mesmo após o processo de abstração, apresentam um impacto positivo no desempenho do agente.

Palavras-chave: Representação Relacional; Políticas Abstratas; Aprendizado por Reforço; Programação Lógica Indutiva; Robótica.

1 INTRODUÇÃO

Atualmente robôs móveis encontram-se presentes nos mais diversos locais, auxiliando em serviços de entrega e limpeza em hospitais, escritórios e residências. A capacidade de reutilização do conhecimento adquirido previamente poderia acelerar o aprendizado de novas tarefas, dessa forma conferindo aos robôs uma maior autonomia e capacidade de adaptação, contribuindo para a aceitação e conseqüente difusão de seu uso. Apesar da diversidade das soluções encontradas na literatura para a navegação autônoma de

robôs móveis [1, 2, 3], elas tipicamente não consideram soluções previamente encontradas para problemas similares. Isso impede que o robô tenha a habilidade de melhorar seu desempenho por meio do aprendizado decorrente de experiências passadas em tarefas similares. Além disso, os planos de navegação gerados por essas soluções poderiam ser mais facilmente comunicados a seres humanos se utilizássemos representações baseadas em Lógica de Primeira Ordem [4].

Tem havido um grande e recente interesse em abordar o problema de adaptação em tarefas de navegação e planejamento em robôs móveis, utilizando uma linguagem mais expressiva. Diversas técnicas têm sido exploradas, desde a construção de descrições relacionais de macroações (planos parciais) para transferência e reuso no aprendizado de novas políticas de atuação [5] até a utilização de aprendizes simbólicos proposicionais para generalizar a política ótima para outros problemas [6].

Este artigo tem por objetivo propor uma alternativa para a transferência do conhecimento obtido na solução prévia de um problema de navegação e aplicação desse conhecimento na solução de um novo problema de navegação através da geração de uma política abstrata com a utilização do algoritmo Qab-Learning. Além disso, visa avaliar comparativamente esta nova técnica com uma proposta anteriormente na literatura [7] [8], em aplicações de navegação robótica em ambientes internos de edifícios.

O artigo está organizado da seguinte forma: a seção II apresenta os conceitos de *Processos Markovianos de Decisão e Lógica de Primeira Ordem*. As seções III e IV descrevem duas técnicas para se aprender políticas abstratas, ND-TILDE e Qab-Learning. A seção V descreve como se dá a aplicação das políticas em um caso de navegação robótica, e os resultados obtidos dos experimentos são discutidos na seção VI. Na seção VII são apresentadas as conclusões e diretrizes para trabalhos futuros.

2 REPRESENTAÇÃO RELACIONAL

Uma abordagem tradicional de formalização para problemas estocásticos de decisão sequencial consiste no uso de Processos Markovianos de Decisão (MDP – *Markov Decision Process*), que podem ser expandidos para MDPs relacionais, ou RMDPs (*Relational Markov Decision Process*), nos quais é possível que as descrições de estados e ações sejam feitas com a utilização de Lógica de Primeira Ordem.

Um MDP pode ser definido formalmente pela quádrupla $\langle S, A, T, R \rangle$ [9], onde S é um conjunto de estados e A de ações, ambos finitos, sendo o conjunto de ações aplicáveis no estado $s \in S$ denotado por $A(s)$. $T: S \times A \times S \rightarrow [0, 1]$ é a *função de transição*, com $T(s_i, a_i, s_{i+1})$ definindo a probabilidade de transição do estado s_i para o estado s_{i+1} através da execução da ação a_i e $R: S \times A \times S \rightarrow \mathcal{R}$ é a *função de recompensa*, com $r_i := R(s_i, a_i, s_{i+1})$. Dessa forma, uma transição do estado $i \in S$ para $j \in S$, decorrente da execução de alguma ação $a \in A(i)$, ocorre com probabilidade $T(i, a, j)$, e então uma recompensa $R(i, a, j)$ é recebida. Em um MDP, pode-se utilizar Lógica de Primeira Ordem na descrição de estados e ações, aplicando variáveis e buscando evidenciar as relações entre os objetos, transformando-o em um RMDP [10, 11]. Expressões da forma $p(t_1, t_2, \dots, t_m)$ são *átomos*, com p sendo um símbolo de relação entre os termos t_i . Por exemplo, $con(r_1, r_2)$ indica que os cômodos r_1 e r_2 estão conectados (são adjacentes e mutuamente acessíveis). Um *termo* pode consistir em uma variável (iniciando com letra maiúscula) ou uma constante (iniciando com letra minúscula, como r_1). Conjunções de átomos lógicos representam estados. Assim, um exemplo de estado $z \in S$ seria:

$$z \equiv in(r_1), con(r_1, r_2), room(r_1), room(r_2).$$

indicando que o robô está no cômodo r_1 , o qual é conectado ao cômodo r_2 .

Da mesma forma, define-se o conceito de *estado abstrato* como uma conjunção de átomos lógicos que contêm variáveis. Por exemplo, o estado abstrato

$$Z \equiv in(R), con(R, R'), room(R), room(R')$$

representa estados nos quais o robô está no cômodo R , o qual é conectado a um outro cômodo R' . Neste caso, o estado Z seria uma instância do estado abstrato Z .

Uma *substituição* θ atribui termos às variáveis de forma que $Z\theta \subseteq z$. Assim, para o exemplo acima, a instância z seria alcançada com a substituição $\theta = \{R/r_1, R'/r_2\}$.

Um estado é denominado *concreto* se não contiver variáveis. Conjunções com variáveis são consideradas existencialmente quantificadas. Uma conjunção Z é *subsumida* pela conjunção W , isto é, $Z \leq_{\theta} W$, se existir θ tal que $W\theta \subseteq Z$.

Analogamente, ações abstratas podem ser representadas por átomos que possuem variáveis como argumentos. Assim,

$goto(R, R')$ é uma ação abstrata que indica que o robô navega de um cômodo R para um cômodo R' , e $goto(r_1, r_2)$ é uma ação concreta com a substituição $\theta = \{R/r_1, R'/r_2\}$.

Resolver um problema modelado por um (R)MDP consiste em computar uma *política* π , que especifica quais ações $a \in A$ devem ser executadas quando o agente se encontra no estado $s \in S$. Se a política for *não determinística*, tem-se que $\pi: S \rightarrow 2^A$; entretanto, se a política for *determinística*, tem-se que $\pi: S \rightarrow A$, ou seja, π é determinística se leva cada estado em S a uma única ação em A .

Dada a política π e um *fator de desconto* $\gamma \in [0, 1]$, que desconta recompensas futuras, a *função de valor de estado* $V^{\pi}: S \rightarrow \mathcal{R}$ representa o valor esperado de recompensas futuras que um agente espera acumular ao seguir a política π .

A *política ótima* π^* é a política que satisfaz a condição: $V^{\pi^*} \geq V^{\pi}$, $\forall s \in S$ e $\forall \pi$, sendo a *função de valor de estado ótima* denotada por V^* . Assim, o interesse consiste em computar uma política que melhor aproxime ou, idealmente, iguale a política ótima π^* .

Da mesma forma que foram definidos estados e ações abstratas no RMDP, pode-se definir uma *política abstrata* π_a [11] como uma lista ordenada de regras de ações abstratas na forma $i: \sigma \Rightarrow \{\alpha\}$, sendo σ um estado abstrato, $\{\alpha\}$ um conjunto de ações abstratas e i um inteiro representando a ordenação da regra.

Existem diversas formas para encontrar a política que soluciona um (R)MDP. Um exemplo é a utilização de algoritmos de planejamento de atividades, como o SPUDD [12], ou técnicas de aprendizado que utilizam repetidas iterações do agente com o ambiente, como no aprendizado por reforço [9]. Por sua vez, uma política abstrata pode ser obtida tanto por meio da generalização de uma política concreta conhecida quanto pelo aprendizado direto. Neste artigo, essas duas formas são exploradas. A primeira utiliza uma política concreta conhecida para gerar exemplos, que são agregados por meio de aprendizado supervisionado usando o algoritmo ND-TILDE [7], descrito na seção III. Já a segunda forma utiliza aprendizado por reforço [9] por meio da abordagem proposta nesse artigo, descrita na seção IV.

3 APRENDENDO POLÍTICAS ABSTRATAS COM APRENDIZADO INDUTIVO.

Esta seção descreve um método de obtenção de uma política abstrata através da indução de uma árvore de decisão em Lógica de Primeira Ordem [8] a partir de exemplos gerados por uma política concreta conhecida.

Dado um conjunto de exemplos E , onde cada exemplo consiste em um par estado-ação obtido durante a execução de uma política concreta conhecida, pode-se aplicar o algoritmo ND-TILDE, descrito em pseudo-código no Algoritmo 1. Ele induz a política abstrata a partir de E , gerando uma

representação através de uma árvore binária de decisão em Lógica de Primeira Ordem.

Algoritmo 1: Algoritmo ND-TILDE

```

1: função ND-TILDE (E: cj. de exemplos): retorna uma
   árvore binária de decisão
2:   T ← GERAR_TESTES(E)
3:   τ ← DIVISAO_OTIMA(T,E)
4:   Ee ← TESTE_VERDADEIRO(E,τ)
5:   Ed ← TESTE_FALSO(E,τ)
6:   se CRITERIO_PARADA(E,Ee,Ed) então
7:     retorna folha(INFO(E))
8:   senão
9:     te ← ND-TILDE(Ee)
10:    td ← ND-TILDE(Ed)
11:    retorna nó_interno(τ,te,td)
12:  fim se
13: fim função
    
```

ND-TILDE utiliza um conjunto de exemplos de treinamento E para criar uma árvore. Os candidatos aos testes de decisão são criados a partir de um conjunto de operadores de refinamento [8] definidos previamente por um especialista (passo 2 do Algoritmo 1). Cada um desses operadores gera um conjunto de literais de primeira ordem que é candidato para a divisão do conjunto de exemplos E fornecido. O melhor teste que pode ser aplicado para dividir o conjunto E é aquele que apresenta a maior redução de entropia. O teste ótimo é escolhido utilizando-se a taxa de ganho padrão. Esse teste particiona o conjunto E de exemplos em dois outros conjuntos: E_e que contém os exemplos onde o teste τ é verdadeiro (passo 4); e E_d que contém os exemplos onde o teste τ é falso (passo 5).

Se a partição induzida em E indica que a divisão deva parar (procedimento CRITERIO_PARADA no passo 6), uma folha contendo sentenças atômicas que representam ações abstratas é então criada na árvore (passo 7). Caso mais de uma sentença atômica representando uma ação esteja associada aos exemplos remanescentes na folha, todas elas são adicionadas pelo algoritmo a esta folha, gerando uma política não-determinística. Se a partição induzida em E não indica que o processo deva parar, então para cada um dos elementos da partição (conjuntos E_e e E_d), a função ND-TILDE é chamada recursivamente (passos 9 e 10). Um nó interno é então criado (passo 11) usando o teste ótimo τ como teste e tendo como filhos as duas sub-árvores (τ_e e τ_d) criadas em cada chamada de ND-TILDE.

Ao término do algoritmo, cada caminho na árvore de decisão gerada da raiz até uma folha (excluindo esta) determina um estado abstrato, e cada folha possui um conjunto de ações abstratas relacionadas àquele estado. Percorrendo essa árvore sistematicamente seguindo a

estratégia de busca em profundidade, obtém-se a sequência de regras ordenadas que compõem a política abstrata. A Seção V apresenta um exemplo de política abstrata para o problema da navegação em robótica obtida através da aplicação do ND-TILDE.

4 APRENDENDO POLÍTICAS ABSTRATAS COM APRENDIZADO POR REFORÇO

Nessa seção é proposto o uso de aprendizado por reforço (RL – *Reinforcement Learning*) como uma alternativa para o aprendizado da política abstrata que generaliza o conhecimento da solução de um problema defrontado pelo robô. No RL, o aprendizado se dá por meio da interação direta do agente com o ambiente, em intervalos de tempos discretos contendo ciclos alternados de percepção e ação.

Tido como o mais popular algoritmo de RL, o algoritmo Q-Learning é um algoritmo de aprendizado livre de modelos, com o qual uma política ótima π* é aprendida iterativamente quando o modelo do sistema não é conhecido, i.e., as funções T e/ou R do (R)MDP são desconhecidas [9]. Fazendo V*(s) = max_a Q*(s,a), o algoritmo Q-Learning iterativamente aproxima a estimativa de Q*(s,a), da seguinte forma:

$$Q_{t+1}(s_t, a_t) \leftarrow (1 - \alpha) Q_t(s_t, a_t) + \alpha (r_t + \gamma \max_a Q_t(s_{t+1}, a)), \quad (1)$$

sendo s_t o estado atual, a_t a ação realizada em s_t, r_t o reforço recebido após realizar a_t em s_t e atingir s_{t+1}, s_{t+1} o novo estado, γ ∈ [0, 1] o fator de desconto e α ∈]0, 1] a taxa de aprendizagem. O aprendizado aqui conduzido segue uma política ε-greedy, que consiste em aplicar uma ação totalmente aleatória quando o valor de uma variável aleatória for menor que ε (correspondendo à fase de exploração), e aplicar a política greedy, isto é, π(s) = arg max_a Q(s, a), no caso contrário (correspondendo à fase de exploração).

Algoritmo 2: Qab-Learning

```

1: função Qab-Learning (σ: cj. de estados abstratos, α: cj. de
   ações abstratas, ε: fator exploração): retorna uma política
   abstrata
2:   INICIALIZAR Qab0
3:   Observar o estado concreto st
4:   Determinar o estado abstrato σt que subsume st, isto é,
      st ≤0t σt
5:   repetir até CONDIÇÃO_PARADA
6:     se um número aleatório ∈ [0,1] ≤ ε então
7:       Determinar a ação concreta at = RANDOM(A(σt))
8:     senão
9:       Determinar ação abstrata: αt ← arg maxα Qab(σt,
      α)
10:    Determinar o cj. de ações concretas subsumidas
      por αt com a substituição θt: {at} ≤0t αt
    
```

```

11:   Fazer  $a_t = \text{RANDOM}(\{a_t\})$ 
12:   fim se
13:   Executar a ação concreta  $a_t$ 
14:   Observar o novo estado concreto  $s_{t+1}$  e o reforço recebido  $r_t$ 
15:   Atualizar  $Q_{ab}(\sigma_t, \alpha_t)$  usando a equação (1), com  $\sigma_t$ ,  $\sigma_{t+1}$  e  $\alpha_t$ , que subsumem  $s_t$ ,  $s_{t+1}$ , e  $a_t$ , respectivamente
16:    $s_t \leftarrow s_{t+1}$ ;  $\sigma_t \leftarrow \sigma_{t+1}$ 
17:   fim repetir
18:   laço para cada  $\sigma \in \sigma$  faça
19:      $\pi_{abs}(\sigma) \leftarrow \arg \max_{\alpha} Q_{ab}(\sigma, \alpha)$ 
20:   fim laço
21:   retornar ( $\pi_{abs}$ )
22: fim função
    
```

O Algoritmo 2 apresenta o pseudo-código do algoritmo Qab-Learning, que é uma modificação do algoritmo Q-Learning para aprender diretamente uma política abstrata. O algoritmo Qab-Learning consiste basicamente no Q-Learning tradicional usando a política ϵ -greedy, porém o aprendizado se dá no nível abstrato. Observa-se o estado concreto atual (passo 3) e determina-se o estado abstrato que subsume esse estado concreto e a correspondente substituição θ_t (passo 4). O ciclo de percepção e ação prossegue até atingir uma condição de parada (passo 5). Se na iteração atual uma exploração deve ser feita (passo 6), então deve-se determinar uma ação aleatória concreta do conjunto de ações aplicáveis no estado concreto atual (passo 7). Se não, fazer exploração determinando a ação abstrata recomendada pela política *greedy* para o estado abstrato que subsume o estado concreto atual (passo 9) e determinar o conjunto de ações concretas subsumidas por esta ação abstrata (passo 10). Escolher aleatoriamente uma ação concreta desse conjunto (passo 11). Executar a ação definida pela fase de exploração ou de exploração (passo 13), observar o reforço recebido e o novo estado concreto (passo 14). Determinar o estado abstrato que subsume o novo estado concreto observado e atualizar a tabela Qab no nível abstrato com esses dados (passo 15). Atualizar os estados concreto e abstrato atuais correspondentes ao estado concreto observado (passo 16) e repetir o ciclo. No final, retornar a política aprendida no nível abstrato. Na próxima seção esse algoritmo é aplicado ao problema de navegação robótica, gerando uma política abstrata que busca generalizar a solução do problema.

5 APLICANDO NA NAVEGAÇÃO ROBÓTICA

Considere agora um problema de navegação de robôs, cujo ambiente é o ilustrado na Fig. 1 (doravante *ambiente 1*).

Os estados desse problema podem ser descritos de acordo com o seguinte conjunto de predicados: $\{\text{corridor}/1, \text{room}/1, \text{center}/1, \text{nearDoor}/1, \text{in}/1, \text{con}/2\}$, e o conjunto de ações com: $\{\text{gotoRDRD}/2, \text{gotoCDCD}/2,$

$\text{gotoCCCC}/2, \text{gotoRCRD}/2, \text{gotoRDRC}/2, \text{gotoRDCD}/2, \text{gotoCDRD}/2, \text{gotoCCCD}/2, \text{gotoCDCC}/2\}$.

O significado de cada predicado é literal: por exemplo, se o agente estiver ocupando a localização 1 no mapa (representado como objeto *l1*), esse estado pode ser descrito como:

$s_{11} = \{\text{in}(l1), \text{corridor}(l1), \text{center}(l1), \text{con}(l1, l5), \text{corridor}(l5), \text{nearDoor}(l5)\}$

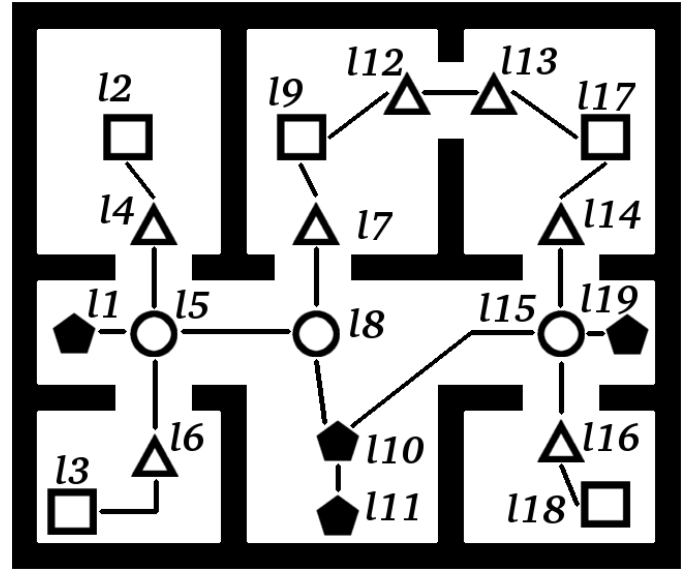


Figura 1 – Mapa do ambiente utilizado para o aprendizado de políticas (ambiente 1). Centros e portas de cômodos são representados por quadrados e triângulos, respectivamente, e círculos e pentágonos pretos representam localidades em corredores próximos e distantes de portas, respectivamente. As linhas que conectam as localidades, duas a duas, na figura, indicam que essas localidades são adjacentes.

Se o agente faz a escolha de ir da localização *l1*, que é *corridor(l1)* e *center(l1)*, para a localização *l5*, que é *corridor(l5)* e *nearDoor(l5)*, descreve-se essa ação como *gotoCCCD(l1, l5)*. As abreviações CC e CD indicam *corridor* e *center*, *corridor* e *door*, respectivamente, e as demais seguem o mesmo padrão.

A pré-condição de *gotoXXYY(Li, Lj)* é $(\text{in}(Li) \wedge \text{con}(Li, Lj))$ e ambas localidades envolvidas na ação são representadas pelas variáveis *Li* e *Lj*, as quais atendem às condições *XX* e *YY*.

A. Definindo as Políticas Abstratas

O problema que o agente deve resolver é episódico e consiste em atingir um estado-meta. Esse estado-meta foi modelado como sendo o único estado para o qual o agente recebe um reforço positivo, independentemente da posição inicial.

Dado esse domínio, a tarefa utilizada para a indução das duas políticas abstratas foi a de ir de qualquer localidade no ambiente 1 para a localidade *l2*. Ou seja, o objetivo é chegar a *l2*, partindo de qualquer outro lugar no mapa. A primeira política, π_{ND} , foi criada com o algoritmo ND-TILDE, fornecido na Seção III (Algoritmo 1), e a segunda, π_{Qab} , com

o algoritmo Qab-Learning, fornecido na Seção IV (Algoritmo 2).

Para a política π_{ND} , em primeiro lugar descobriu-se a política ótima concreta para o ambiente 1. Dada essa política conhecida, aplicou-se o ND-TILDE para a extração da política abstrata π_{ND} em forma de árvore de decisão. A árvore pode também ser descrita como um conjunto de regras ordenadas, que estão ilustradas na Tabela I.

Nota-se que existem 8 estados abstratos nessa política, correspondentes às partições que o algoritmo julgou serem as mais relevantes para o problema. Além disso, cada estado aponta para um conjunto de ações abstratas. Na implementação dos experimentos deste artigo, descritos na Seção VI, a escolha entre uma das ações deste conjunto é feita aleatoriamente.

Para a execução de π_{Qab} , cada descrição de estado concreto (a descrição de um estado só contempla conexões com localidades vizinhas) foi transformada em um estado abstrato simplesmente substituindo por variável cada constante presente nos termos dos predicados, resultando em 12 estados abstratos diferentes, que generalizam os 19 estados concretos do problema. De posse desse conjunto de estados abstratos e do conjunto de ações abstratas, o algoritmo Qab-Learning (Algoritmo 2) foi usado para gerar a política abstrata π_{Qab} , presente na Tabela II.

TABELA I
POLÍTICA ABSTRATA CONSTRUÍDA COM O ND-TILDE

i	Estado Abstrato	Ações Abstratas
1	$in(A), con(A,B), center(B),$ $corridor(B), center(A),$ $corridor(A), con(A,C),$ $nearDoor(C), corridor(C)$	$\{gotoCCDC(A,C)\}$
2	$in(A), con(A,B), center(B),$ $corridor(B), center(A), corridor(A)$	$\{gotoCCCC(A,B)\}$
3	$in(A), con(A,B), center(B),$ $corridor(B), con(A,C),$ $nearDoor(C), corridor(C)$	$\{gotoCDRD(A,D),$ $gotoCDCD(A,C)\}$
4	$in(A), con(A,B), center(B),$ $corridor(B)$	$\{gotoCDCC(A,B)\}$
5	$in(A), nearDoor(C), corridor(C),$ $center(A), corridor(A)$	$\{gotoCCCD(A,C)\}$
6	$in(A), nearDoor(C), corridor(C)$	$\{gotoRDRC(A,D)$ $gotoRDCA(A,C)\}$
7	$in(A), room(A), nearDoor(A)$	$\{gotoRDRC(A,D),$ $gotoRDRD(A,D)\}$
8	$in(A)$	$\{gotoRCRD(A,D)\}$

TABELA II
POLÍTICA ABSTRATA CONSTRUÍDA COM QAB-LEARNING

i	Estado Abstrato	Ações Abstratas
1	$in(A), con(A,B), con(A,C),$ $con(A,D), con(A,E),$ $nearDoor(A), corridor(A),$ $center(B), corridor(B),$ $nearDoor(C), room(C),$	$gotoCDRD(A,D)$

	$nearDoor(D), room(D),$ $nearDoor(E), corridor(E)$	
2	$in(A), con(A,B), con(A,C),$ $con(A,D), con(A,E), nearDoor(),$ $corridor(A), center(B),$ $corridor(B), nearDoor(),$ $room(C), nearDoor(D), room(D)$ $, center(E), corridor(E)$	$gotoCDCC(A,B)$
3	$in(A), con(A,B), con(A,C),$ $con(A,D), nearDoor(A),$ $corridor(A), nearDoor(B),$ $corridor(B), nearDoor(C),$ $room(C), center(D), corridor(D)$	$gotoCDCA(A,B)$
4	$in(A), con(A,B), con(A,C),$ $con(A,D), center(A), corridor(A),$ $nearDoor(B), corridor(B),$ $center(C), corridor(C),$ $nearDoor(D), corridor(D)$	$gotoCCCD(A,B)$
5	$in(A), con(A,B), nearDoor(A),$ $room(A), center(B), room(B),$ $con(A,C), nearDoor(C),$ $room(C), con(A,D),$ $nearDoor(D), room(D)$	$gotoRDRD(A,C)$
6	$in(A), con(A,B), nearDoor(A),$ $room(A), center(B), room(B),$ $con(A,C), nearDoor(C),$ $room(C), con(A,D),$ $nearDoor(D), corridor(D)$	$gotoRDCA(A,D)$
7	$in(A), con(A,B), nearDoor(A),$ $room(A), center(B), room(B),$ $con(A,C), nearDoor(C),$ $corridor(C)$	$gotoRDRC(A,D)$ ou $gotoRDCA(A,D)$
8	$in(A), con(A,B), con(A,C),$ $center(A), room(A),$ $nearDoor(B), room(B),$ $nearDoor(C), room(C)$	$gotoRCRD(A,B)$
9	$in(A), con(A,B), con(A,C),$ $nearDoor(A), room(A),$ $center(B), room(B),$ $nearDoor(C), room(C)$	$gotoRDRC(A,B)$
10	$in(A), con(A,B), center(A),$ $corridor(A), nearDoor(B),$ $corridor(B)$	$gotoCCDC(A,B)$
11	$in(A), con(A,B), center(A),$ $room(A), nearDoor(B), room(B)$	$gotoRCRD(A,B)$
12	$in(A), con(A,B), center(A),$ $corridor(A), center(B),$ $corridor(B)$	$gotoCCCC(A,B)$

B. Aplicando as Políticas Abstratas

O objetivo da construção dessas políticas abstratas é possibilitar seus usos em outros problemas similares, em que a extensão da transferência desse conhecimento adquirido pode ser avaliada. Para essa aplicação em robótica móvel, as políticas geradas foram testadas no mesmo ambiente 1, porém

para tarefas distintas, e também em um outro ambiente (ambiente 2), maior que o anterior e cujo mapa pode ser observado na Fig. 2. Para a escolha deste problema “similar”, foram consideradas duas características para serem mantidas fixas entre os problemas: o conjunto de predicados que descrevem os estados e o conjunto de ações abstratas.

A aplicação da política abstrata em um problema ocorre da seguinte forma: primeiro, observa-se o estado concreto s_t no qual o agente se encontra. Busca-se nas regras que definem a política abstrata em uso (π_{ND} ou π_{Qab}), na ordem, qual estado abstrato σ_t (antecedente da regra) que subsume o estado s_t , isso é, $s_t \leq_{\theta} \sigma_t$, definindo a regra e a respectiva substituição θ . Dessa forma obtém-se o correspondente conjunto de ações abstratas com a respectiva aplicação da substituição θ . De posse desse conjunto de ações abstratas instanciadas com a substituição θ , uma delas é escolhida aleatoriamente (chamada aqui de α_{θ}). Caso essa ação α_{θ} escolhida já seja uma ação concreta (devido à substituição aplicada), essa ação concreta a_t é executada em s_t ; caso contrário, ela definirá um conjunto $A\alpha_{\theta}$ de ações concretas que poderão ser aplicadas em s_t (substituindo adequadamente todas as variáveis restantes em α_{θ}), com $A\alpha_{\theta} \subseteq A(s_t)$, e uma escolha aleatória é feita nesse conjunto $A\alpha_{\theta}$ para definir a ação concreta a_t a ser executada em s_t . Para o caso de aplicar-se em um outro ambiente a política abstrata gerada para um determinado ambiente, é possível que o estado concreto atual s_t não seja subsumido por nenhum estado abstrato da política abstrata em uso; neste caso, a implementação corrente escolhe aleatoriamente uma ação a_t no conjunto $A(s_t)$ e executa a_t em s_t . Esse processo é então repetido para o novo estado s_{t+1} , alcançado com a aplicação de a_t em s_t , até que s_{t+1} seja um estado-meta.

Para acelerar a execução da tarefa, quando o agente se depara com um estado no qual ele já passou anteriormente no mesmo episódio, todas as escolhas aleatórias citadas são feitas excluindo-se a(s) escolha(s) anterior(es). Ainda, caso o agente volte a um mesmo estado após ter tentado todas as ações possíveis indicadas pela sua política, então a política aleatória é empregada, para que o agente possa assim encontrar sua meta.

Para exemplificar este procedimento, considere que o ambiente seja o ambiente 1 (Fig.1), que esteja sendo aplicada a política abstrata π_{ND} (Tabela I) e que o agente se encontre no local $l8$, sendo

$$s_t = \{in(l8), con(l8,l10), center(l10), corredor(l10), con(l8, l5), nearDoor(l5), corredor(l5), con(l8, l7), nearDoor(l7), room(l7)\}.$$

Aplicando a política π_{ND} , verifica-se que o estado abstrato que subsume s_t é o estado da terceira regra da Tabela 1:

$$\sigma_t = \{in(A), con(A,B), center(B), corredor(B), con(A,C), nearDoor(C), corredor(C)\},$$

com substituição $\theta = \{A/l8, B/l10, C/l5\}$, e o conjunto de ações abstratas definido pela política π_{ND} é $\{gotoCDRD(A,D), gotoCDCD(A,C)\}$. Supondo que a escolha aleatória recaia

sobre a ação abstrata $gotoCDRD(A,D)$, a aplicação de θ resulta em

$$A\alpha_{\theta} = \{gotoCDRD(l8,D)\} = \{gotoCDRD(l8,l7)\},$$

já que $D=l7$ é a única substituição possível nesse caso. Assim, a ação concreta $a_t = gotoCDRD(l8,l7)$ é executada em s_t , e o agente então observa o novo estado e repete o procedimento, até atingir a meta.

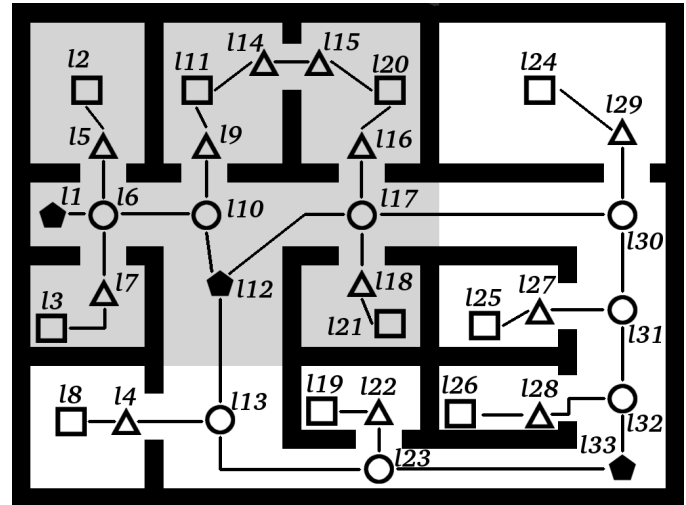


Figura 2 - Mapa do ambiente utilizado para a transferência de políticas (ambiente 2). Centros e portas de cômodos são representados por quadrados e triângulos, respectivamente, e círculos e pentágonos pretos representam localidades em corredores próximos e distantes de portas, respectivamente. As linhas que conectam as localidades, duas a duas, na figura, indicam que essas localidades são adjacentes. A região marcada, similar ao ambiente 1, foi utilizada nos experimentos para verificar como a presença, no novo problema, de um subproblema já resolvido, impactava no desempenho do aprendizado.

6 EXPERIMENTOS E ANÁLISE COMPARATIVA

Para avaliar as políticas abstratas aprendidas, uma série de experimentos foi realizada, usando tanto o ambiente 1, onde as políticas abstratas π_{ND} e π_{Qab} foram geradas, quanto o ambiente 2.

As três políticas utilizadas nos testes são: (i) uma política totalmente aleatória, que serve como referência, (ii) a política abstrata π_{ND} gerada através de aprendizado supervisionado utilizando ND-TILDE (Algoritmo 1) e (iii) a política abstrata π_{Qab} obtida através de aprendizado por reforço utilizando o Qab-Learning (Algoritmo 2). Os testes foram realizados com o objetivo de avaliar o desempenho das políticas na resolução de diversas tarefas.

Um episódio corresponde a iniciar o robô em um local qualquer e aplicar a política que está sendo avaliada até que ele atinja a meta, caso em que temos um episódio bem-sucedido, ou até que um limite máximo de 1.000 passos seja atingido sem que o agente tenha alcançado a meta, e nesse caso considera-se que o episódio não teve sucesso. A cada experimento está associado um conjunto de metas e para cada uma delas foram executados 10.000 episódios. Os gráficos presentes neste artigo consideram a média dos valores obtidos na execução de cada episódio.

O eixo das abscissas em cada gráfico representa o número de passos n , enquanto o eixo das ordenadas representa a fração acumulada de episódios, dentre os 10.000, nos quais o robô atingiu a meta (em outras palavras, se definirmos que N é o número de passos usados para atingir a meta, essa fração representa a distribuição acumulada de N , i.e., $P(N \leq n)$). Ou seja, supondo $n=100$ passos no gráfico da Fig. 3, tem-se que, usando a política aleatória, uma fração de cerca de 60% dos episódios alcançaram a meta em até 100 passos, dentre os episódios que tiveram sucesso (isto é, atingiram a meta em 1.000 passos ou menos) dos 10.000 episódios executados; usando a política π_{ND} e π_{Qab} , as frações foram aproximadamente 78% e 80%, respectivamente.

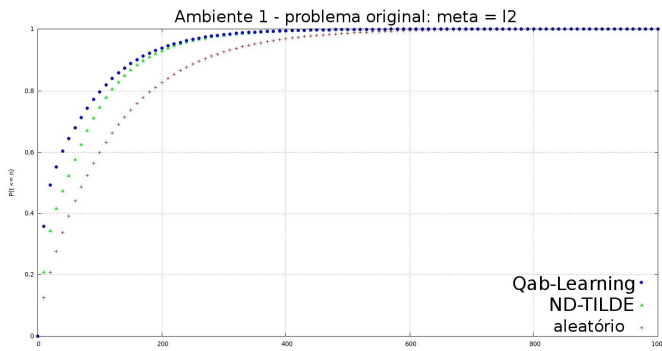


Figura 3 – Resultados da aplicação das políticas abstratas e aleatória no problema original: no ambiente 1, sair de qualquer lugar (exceto o local $l2$) e chegar a $l2$.

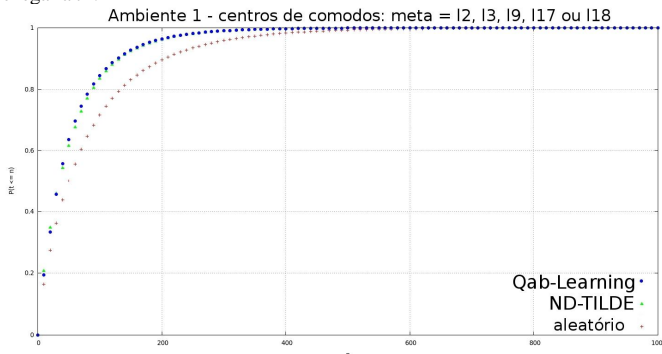


Figura 4 – Resultados da aplicação das políticas abstratas e aleatória no ambiente 1 para a tarefa de, alcançar cada um dos centros de cômodos do ambiente 1 a partir das demais localizações que não a que se pretende alcançar. O gráfico apresenta a média dos resultados obtidos tendo como metas em $l2$, $l3$, $l9$, $l17$ e $l18$, em que cada experimento relacionado a cada uma das metas consistiu em 1000 episódios.

O primeiro experimento consistiu em apenas aplicar as políticas abstratas (e aleatória) para resolver o problema para o qual elas foram geradas, ou seja, o mesmo ambiente (ambiente 1) com a mesma meta (alcançar o local $l2$). Isso foi feito para verificar os efeitos e perdas causados pela abstração. Os resultados estão ilustrados na Fig. 3, onde se observa que o Qab-Learning mostrou melhor desempenho nesse caso. Ainda no mesmo ambiente, outro experimento avaliou a execução das políticas na tarefa de, saindo de qualquer local no ambiente 1, atingir um dos centros de cômodo (locais $l2$, $l3$, $l9$, $l17$ ou $l18$), apresentando a média alcançada pela execução de todas estas tarefas (cada uma

consistindo em 10.000 episódios com no máximo 1000 passos cada). O resultado deste segundo experimento está ilustrado na Fig. 4, onde se verifica que o Qab-Learning e ND-TILDE apresentam resultados similares e melhores do que a política aleatória.

Outros três experimentos foram realizados no mapa da Fig. 2 com o objetivo de avaliar o desempenho das políticas abstratas definidas para o ambiente 1 quando aplicadas no ambiente 2. Nesses experimentos, buscou-se avaliar a capacidade de transferência do conhecimento adquirido em um problema para um novo problema. O experimento, cujo resultado está ilustrado na Fig. 5, mostra o desempenho das políticas abstratas e aleatória no ambiente 2, quando a tarefa é atingir $l2$ partindo de cada uma das demais localizações. A Fig. 6 mostra o desempenho das políticas em uma tarefa similar, na qual a meta é $l25$ a partir das demais localizações do ambiente 2. Já a Fig. 7 mostra o desempenho apresentado nas tarefas de, saindo de qualquer localização a menos da meta, atingir cada um dos centros das salas do ambiente 2.

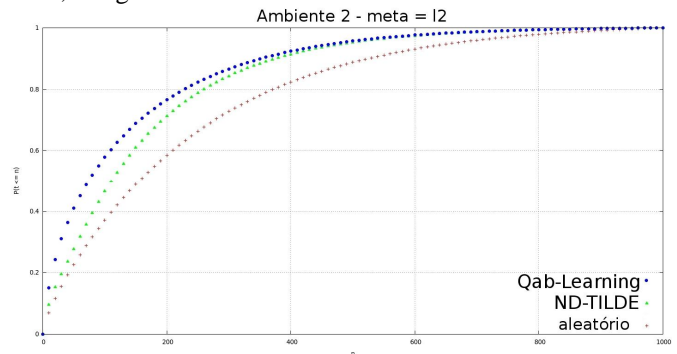


Figura 5 – Resultados da aplicação das políticas abstratas e aleatória no ambiente 2, na tarefa de, saindo de qualquer uma das demais localidades, chegar a $l2$.

Analisando as figuras, pode-se observar que, tanto ND-TILDE quanto Qab-Learning geram políticas que apresentam um melhor desempenho que o resultado pelo uso da política aleatória.

O desempenho superior da política π_{Qab} em relação à π_{ND} para a solução do problema que foi utilizado na obtenção de ambas (Fig. 3), sugere que π_{Qab} foi capaz de reter mais elementos necessários à solução do problema original do que π_{ND} , isto é, a política π_{Qab} seria mais *especializada* que π_{ND} . Entretanto, quando a tarefa é mais geral (Fig. 4), ou seja, atingir uma localidade que, apesar de possuir a mesma natureza (centro de cômodo), possui uma localização diferente e, conseqüentemente, uma caracterização topológica dada por sua vizinhança também diferente (Fig. 1), ambas as políticas apresentam um desempenho similar, superando a aleatória.

Esse resultado, aliado ao que foi discutido anteriormente, poderia apoiar uma interpretação positiva em relação à π_{Qab} , sugerindo que essa política, além de reter um relativo bom desempenho na solução do problema que foi treinada para resolver, possui também características que poderiam ser explicadas por uma capacidade de generalização, também

apresentada pela π_{ND} . Isso lhes permitiria atingir um bom desempenho em uma tarefa distinta da presente em seus treinos, mas que pertence ao mesmo domínio.

Os resultados que compõem o cerne da proposta desse trabalho são ilustrados nas figuras 5, 6 e 7 e se referem à aplicação das políticas aprendidas no ambiente 1 na execução de tarefas no ambiente 2. No experimento cujos resultados são ilustrados na Fig. 5, o agente parte de cada uma das localizações e deve atingir a meta $l2$ executando, no máximo, 1000 passos. É interessante notar que a localidade $l2$ se encontra em uma região similar ao ambiente 1, correspondendo à meta “antiga”. Pode-se ver que o desempenho das políticas é superior ao da aleatória, sendo que π_{Qab} , novamente, apresenta desempenho superior ao de π_{ND} . Esse fato corrobora a hipótese discutida acima de retenção de elementos do problema utilizado na obtenção da política, isto é, de maior *especialização* de π_{Qab} em relação à π_{ND} .

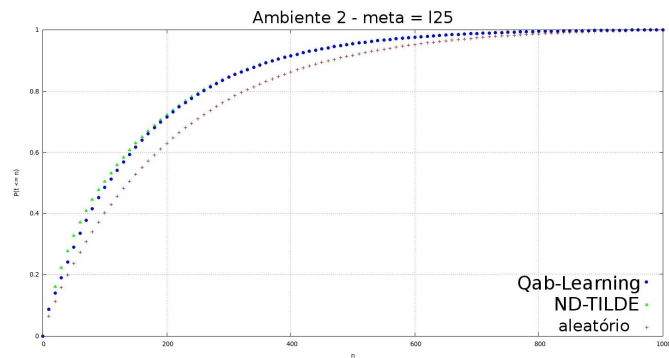


Figura 6 – Resultados da aplicação das políticas abstratas e aleatória no ambiente 2, na tarefa de sair de qualquer lugar e chegar ao local $l25$.

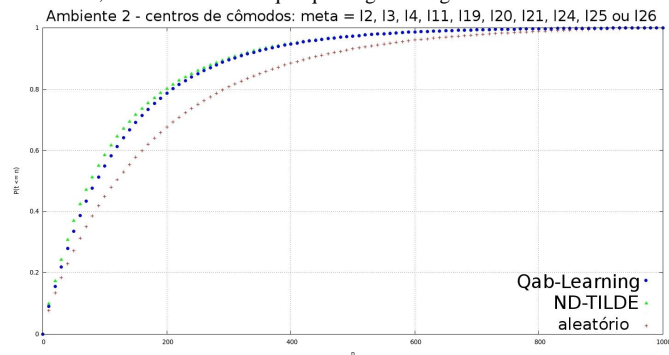


Figura 7 – Resultados da aplicação das políticas abstratas e aleatória no ambiente 2, na tarefa de sair de qualquer lugar e chegar a cada um dos 10 centros dos cômodos. O resultado apresentado é a média dos desempenhos apresentados em cada uma das 10 tarefas.

No experimento cujos resultados são representados na Fig. 6, o agente deveria atingir $l25$, que além de estar fora da região *conhecida* correspondente ao problema original, localiza-se em uma região oposta à $l2$, utilizada como meta na obtenção das políticas. Esperava-se com esse experimento verificar a capacidade de aplicação das políticas quando o problema de navegação consiste em atingir uma localidade que, apesar de possuir a mesma natureza (centro de cômodo), e com isso assemelhar-se à solução do problema original,

possuísse uma caracterização topológica que dificultasse a aplicação direta da política na solução do problema (atingir uma meta em região oposta e com vizinhança ainda desconhecida).

Nesse experimento, as políticas apresentaram um desempenho similar; entretanto, a π_{ND} foi ligeiramente melhor. Uma possível explicação para esse fato poderia ser a de que a suposta retenção de elementos do problema original por π_{Qab} seria acompanhada de *especialização*. Apesar de a diferença entre os desempenhos das políticas ser relativamente pequena, ela sugere que poderia ser positivo analisar o impacto das características intrínsecas ao processo de obtenção das políticas, além da representação utilizada em cada uma, para que se possa inferir quais elementos interferem (e com que magnitude) na capacidade de generalização de uma política concreta obtida na solução de um problema prévio. Uma possível solução para o problema de especialização seria desenvolver e aplicar representações que permitissem a generalização da política abstrata para novos problemas, entretanto sem perder a característica desejável de alto desempenho na solução de um problema já resolvido.

Essa análise é respaldada pelos resultados apresentados na Fig. 7, em que o agente deve atingir os centros de cômodos presentes no ambiente 2 a partir de cada uma das demais localidades. A política gerada pelo ND-TILDE apresenta um desempenho ligeiramente superior, e ambas superam a política abstrata na solução do problema.

7 CONCLUSÕES

Neste artigo, foi proposto o algoritmo Qab-Learning como uma alternativa para a transferência de conhecimento entre tarefas no domínio de navegação robótica, onde o conhecimento da solução prévia de um problema de navegação foi aplicado na solução de um novo problema através da geração de uma política abstrata.

A política abstrata gerada pelo Qab-Learning, através de uma modificação do algoritmo Q-learning, em que estados e ações abstratos são utilizados no aprendizado, é comparada à política gerada pelo algoritmo ND-TILDE. Ambas as políticas geradas por estes algoritmos apresentam um desempenho superior ao de uma política aleatória, i.e., uma escolha aleatória de ações para guiar o robô pelo ambiente. Isso pode ser visto como transferência do aprendizado entre problemas através da utilização das políticas abstratas geradas pelos algoritmos Qab-Learning e ND-TILDE.

Nos experimentos realizados, foi possível verificar que o algoritmo proposto, Qab-Learning, criou políticas abstratas que obtiveram desempenho superior ao da política criada com ND-TILDE em determinadas situações, mas, que ao mesmo tempo, apresentavam desempenho similar ou inferior em outras. Em vista disto, futuramente pretende-se analisar os elementos presentes na obtenção de ambas, buscando verificar quais deles possuem maior e menor impacto na geração de políticas abstratas.

Entretanto, nas situações de superioridade do Qab-Learning, a diferença entre as políticas é notável, o que poderia ser explicado pelo fato de que os processos de geração de estados abstratos através da agregação de estados concretos ser realizada de formas distintas pelos algoritmos ND-TILDE e para a execução do Qab-Learning. Certamente, o intuito da construção dessas políticas abstratas não é utilizá-las como únicos guias do agente, mas sim para acelerar o aprendizado de novas tarefas através da transferência do conhecimento prévio. Desse modo, um dos aspectos a se trabalhar no futuro seria a identificação de subtarefas ou partições do espaço de estados nas quais a política abstrata é apropriada para um novo problema. Ou seja, definir quais partes da política apresentam conhecimento relevante e quando aplicá-las.

Além disso, outro aspecto a ser estudado é a definição e avaliação dos predicados que compõem a descrição do domínio, pois eles possuem um papel decisivo na obtenção das políticas abstratas, dado que o método utilizado para defini-los pode ser determinante na qualidade final das políticas geradas durante o processo de aprendizado.

REFERÊNCIAS

- [1] Dudek, G.; Jenkin, M. **Computational Principles of Mobile Robotics**. New York, Cambridge University Press, 2000.
- [2] Choset, H. M.; Lynch, K. M.; Hutchinson, S.; Kantor, G.; Burgard, W.; Kavraki, L. E.; Thrun, S. **Principles of robot motion: theory, algorithms, and implementation**. New York, MIT Press, 2005.
- [3] Thrun, S.; Burgard, W.; Fox, D. *Probabilistic Robotics*. New York, MIT Press, 2005.
- [4] Russell, S. J.; Norvig, P. **Artificial Intelligence: A Modern Approach**. Upper Saddle River, NJ, Pearson Education, Inc., 2nd. edition, 2003.
- [5] Torrey, L.; Shavlik, J.; Waker, T.; Maclin, R. Relational macros for transfer in reinforcement learning. 2008. In: ILP'07 – **Proc. Of the 17th Int. Conf. on Inductive Logic Programming**, 2008.
- [6] Madden, M. G.; Howley, T. Transfer of experience between reinforcement learning environments with progressive difficulty. **Artificial Intelligence Review**, Kluwer Academic Publishers, Norwell, MA, USA, v. 21, p. 375-398, June 2004.
- [7] Matos, T.; Bergamo, Y.; Silva, V. F. da; Costa, A. H. R. Stochastic Abstract Policies for Knowledge Transfer in Robotic Navigation Tasks. In: MICAI'2011 – **Mexican International Conference on Artificial Intelligence**, 2011.
- [8] Blockeel, H.; Raedt, L.D. Top-down induction of first-order logical decision trees. **Artificial Intelligence**, Volume 101, Issues 1-2, p. 285-297. May 1998.
- [9] Sutton, R.S., Barto, A.G.: **Reinforcement Learning: An Introduction**. MIT Press, Cambridge, MA, USA. 1998.
- [10] Kersting, K.; Otterlo, M.V.; and Raedt, L.D. 2004. Bellman goes relational. In **Proc. of 21th Int. Conf. on Machine Learning**, 465–472.
- [11] van Otterlo, M. **The Logic of Adaptive Behaviour**, IOS Press, Amsterdam, 2009.
- [12] Hoey, J., St-Aubin, R., Hu, A.J., Boutilier, C.: SPUDD: Stochastic planning using decision diagrams. In: UAI'99 – **Proc. of Uncertainty in Artificial Intelligence**, Stockholm, Sweden. 1999.

AGRADECIMENTOS

Esta pesquisa obteve o apoio da FAPESP (08/03995-5, 09/04489-9, 09/14650-1) e do CNPq (305512/2008-0).

Rafael agradece a Guilherme Renoldi pelo auxílio nas etapas de depuração do simulador e pelas ideias surgidas durante as construtivas discussões.