

# Lifelong Learning for Disturbance Rejection on Mobile Robots

David Isele  
Univ. of Pennsylvania  
isele@seas.upenn.edu

José Marcio Luna  
Univ. of Pennsylvania  
joseluna@seas.upenn.edu

Eric Eaton  
Univ. of Pennsylvania  
eeaton@seas.upenn.edu

Gabriel V. de la Cruz  
Washington State Univ.  
gabriel.delacruz@wsu.edu

James Irwin  
Washington State Univ.  
james.irwin@wsu.edu

Brandon Kallaher  
Washington State Univ.  
brandon.kallaher@wsu.edu

Matthew E. Taylor  
Washington State Univ.  
taylorm@eecs.wsu.edu

## ABSTRACT

Learning controllers for multiple systems is often an expensive process when controllers for each system are learned individually. Advances in lifelong learning suggest that information between systems can be shared, improving the quality of the controllers that are learned. However these results have been largely theoretical, with applications limited to benchmark problems with known dynamics. We show that these methods can be extended to robotic platforms. Particularly we validate our assumptions for transfer learning between tasks with unknown dynamics in order to carry out a disturbance rejection problem. We view this as early work leading up to learning robust fault tolerant control.

## Categories and Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning; I.2.9 [Artificial Intelligence]: Robotics

## General Terms

algorithms, experimentation

## Keywords

lifelong learning, policy gradients, reinforcement learning, robotics

## 1. INTRODUCTION

Our goal is to learn fault tolerant control in multi-agent systems. In order to approach this problem we begin by developing a method that can accomplish disturbance rejection across multiple robots.

In control systems, a perfect model of the dynamics of the system is often necessary to guarantee the stabilization of the system. Policy search approaches have been proposed to deal with the design of controllers in model-free applications. However, given a collection of robots each with a different unknown disturbance, learning a unique control policy for each robot in the system can be very costly. One approach

to reduce the amount of learning is to share information between robots.

Lifelong learning [19] is a promising approach for accomplishing information sharing between different robots. It works on-line, allowing the different systems to be encountered consecutively so *a priori* knowledge of all the different systems is not required at the start of training. It also preserves and possibly improves the models encountered early on, in contrast to transfer methods which only optimize performance on the target system.

However most recent work in lifelong learning [19, 5, 6] has been theoretical, using simulations of benchmark problems with known dynamics to demonstrate knowledge sharing. The contribution of this work is to present the first results of lifelong learning on robotic systems. We do this by applying the PG-ELLA framework [5] to the problem of disturbance rejection on a set of Turtlebots in the Gazebo Simulator. In this work in progress, the obtained results motivate future experiments on real robots.

After discussing related work in section 2, we review policy gradients and describe the particular base-learner we are using for lifelong learning in section 3. We then outline our experiments in section 4, showing how lifelong learning can be used to accomplish disturbance rejection on robotic systems.

## 2. RELATED WORK

Some of the typical problems in control systems are regulation, trajectory tracking, disturbance rejection and robustness [9, 15, 17]. The disturbance rejection problem consists of implementing a controller that allows the plant to fulfill the desired task while compensating for a disturbance that modifies its nominal dynamics. As long as there is an accurate model of the plant, several mathematical artifacts have been provided to handle constant, constant and unknown, time-varying and even stochastic disturbances [8, 9, 15].

However, things get more complicated when the goal is to design a controller that stabilizes a system whose model is not available due to complex internal iterations, uncertainty in the system, event-based dynamics and technological limitations. Unlike typical control system theory where policies are generated from a model, we generate policies based on sampled trajectories generated by simulation using reinforcement learning (RL). Reinforcement learning [11] is often utilized to learn controllers in a model-free settings. Amongst reinforcement learning algorithms, policy gradient

**Appears in:** *Proceedings of the 15th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2016)*, John Thangarajah, Karl Tuyls, Stacy Marsella, Catholijn Jonker (eds.), May 9–13, 2016, Singapore. Copyright © 2016, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

(PG) methods [20, 23] are popular in robotic applications [13, 18] since they accommodate continuous state/action spaces and can scale well to high dimensional spaces. Unlike these methods where the goal is to learn a single controller, we use policy gradients in a lifelong learning setting.

It has been shown that PG methods can be used in a lifelong learning setting [5, 6], however these works focus largely on theory, using benchmark simulations to demonstrate their results. While there are examples of lifelong learning on robots, they tend to focus in skill refinement on a single robot [10, 21] rather than sharing information across multiple robots as we do in our work.

### 3. BACKGROUND

Our approach works by sharing knowledge between different robotic systems. The policy for each system is learned by reinforcement learning. In this section we cover the mathematical framework that supports our experiments on lifelong learning.

#### 3.1 Reinforcement Learning

A reinforcement learning (RL) agent must select sequential actions to maximize its expected return. RL approaches do not require previous knowledge of the system dynamics, instead, the control policies are learned through the interactions with the system. RL problems are typically formalized as Markov Decision Processes (MDPs) with the form  $\langle \mathcal{X}, \mathcal{A}, P, R, \gamma \rangle$  where  $\mathcal{X} \subset \mathbb{R}^{d_x}$  is the set of states,  $\mathcal{A}$  is the set of actions,  $P : \mathcal{X} \times \mathcal{A} \times \mathcal{X} \rightarrow [0, 1]$  is the state transition probability describing the systems dynamics,  $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function and  $\gamma \in [0, 1]$  is the reward discount factor. At each time step  $h$ , the agent is in the state  $\mathbf{x}_h \in \mathcal{X}$  and must choose an action  $\mathbf{a}_h \in \mathcal{A}$  so that it transitions to a new state  $\mathbf{x}_{h+1}$  with state transition probability  $P(\mathbf{x}_{h+1} | \mathbf{x}_h, \mathbf{a}_h)$ , yielding a reward  $r_h$  according to  $R$ . The action is selected according to a policy  $\pi : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$  which specifies a probability distribution over actions given the current state. The goal of RL is to find the optimal policy  $\pi^*$  that maximizes the expected reward.

We use a class of RL algorithms known as Policy Gradient (PG) methods [20], which are particularly well suited for solving high dimensional problems with continuous state and action spaces, such as robotic control [18].

The goal of PG is to use gradient steps to optimize the expected average return:

$$\mathcal{J}(\theta) = \int_{\mathbb{T}} p_{\theta}(\tau) \mathcal{R}(\tau) d\tau, \quad (1)$$

where  $\mathbb{T}$  is the set of all trajectories and  $\mathcal{R}(\tau)$  is the average per-step reward, specifically:

$$\begin{aligned} p_{\theta} &= \prod_{h=0}^H p(\mathbf{x}_{h+1} | \mathbf{x}_h, \mathbf{a}_h) \pi(\mathbf{a}_h, \mathbf{x}_h) , \\ \mathcal{R}(\tau) &= \frac{1}{H} \sum_{h=0}^H r(\mathbf{s}_h, \mathbf{a}_h, \mathbf{s}_{h+1}) . \end{aligned}$$

Most PG methods (e.g. episodic REINFORCE [23], Natural Actor Critic [18], and PoWER [13]) optimize the policy by maximizing a lower bound on the return, comparing trajectories generated by different candidate policies  $\pi$ . In this particular application, the PG method we use in our ex-

periments is finite differences (FD) [12] which optimizes the return directly.

#### 3.2 Finite Differences for Policy Search

The local optimization around an existing policy  $\pi$  parameterized by a parameter matrix  $\theta$  is carried out by computing changes in the policy parameters  $\Delta\theta$  that will increase the expected reward, thus producing the iterative update,

$$\theta_{m+1} = \theta_m + \Delta\theta_m.$$

Gradient-based methods for policy updates follow the gradient of the expected return  $\mathcal{J}$  given a step-size  $\delta$ ,

$$\theta_{m+1} = \theta_m + \delta \nabla_{\theta} \mathcal{J}.$$

In FD gradients, we have a set of  $n$  perturbed policy parameters which are used to estimate the effect of a change in policy parameters:

$$\Delta \hat{\mathcal{J}}_{\mathbf{p}} \approx \mathcal{J}(\theta_m + \Delta\theta_{\mathbf{p}}) - \mathcal{J}_{ref},$$

where  $\Delta\theta_{\mathbf{p}}$  are the individual perturbations for  $\mathbf{p} = [1, \dots, n]$ ,  $\Delta \hat{\mathcal{J}}_{\mathbf{p}}$  is the estimate of their effect on the return, and the  $\mathcal{J}_{ref}$  is a reference return which is usually taken as the return of the unperturbed parameters. By using linear regression we get an approximation of the gradient,

$$\nabla_{\theta} \mathcal{J} \approx (\Delta\theta^{\top} \Delta\theta)^{-1} \Delta\theta^{\top} \Delta \hat{\mathcal{J}}_{\mathbf{p}},$$

where  $\Delta \hat{\mathcal{J}}_{\mathbf{p}}$  contains all the stacked samples of  $\Delta \hat{\mathcal{J}}_{\mathbf{p}}$  and  $\Delta\theta$  contains the stacked perturbations  $\Delta\theta_{\mathbf{p}}$ . This approach is sensitive to the type and magnitude of the perturbations, as well as to the step size  $\delta$ . Since the number of perturbations needs to be as large as the number of parameters, this method is considered to be noisy and inefficient for problems with large sets of parameters.

In our experiments, the policy is represented as a function defined over a parameter matrix  $\theta \in \mathbb{R}^{d_{\theta}}$ . Our goal is to optimize the expected average return of Equation 1.

In order to share information across different learned policies, we incorporate the PG learning process into a lifelong machine learning setting.

#### 3.3 Lifelong Machine Learning

Lifelong learning focuses on learning a set of tasks consecutively while performing well across all tasks. Given a round  $t = 1, \dots, T$  a task  $Z^{(t)}$  is observed. In our setting, each task corresponds to a reinforcement learning problem for an individual robot. We assume that the model associated to  $Z^{(t)}$  is parameterized by a parameter  $\theta^{(t)} \in \mathbb{R}^{d_{\theta_t}}$ . The ideal goal is that prior knowledge about tasks  $Z^{(1)}, \dots, Z^{(t-1)}$  should provide enough information so that the lifelong learning algorithm performs better and faster on  $Z^{(t)}$  while being able to scale as the number of tasks increases.

Following work in both multi-task [14] and lifelong learning [19], we assume there is a shared basis  $\mathbf{L} \in \mathbb{R}^{d_x \times k}$  and a sparse weight vector  $\mathbf{s}^{(t)} \in \mathbb{R}$ , so that the policy parameters  $\theta^{(t)}$  are given by,

$$\theta^{(t)} = \mathbf{L} \mathbf{s}^{(t)}.$$

Using the return function in (1) we propose the following multi-task objective function,

$$\argmin_{\mathbf{L}, \mathbf{S}} \frac{1}{T} \sum_t \left[ -\mathcal{J}(\theta^{(t)}) + \lambda \|\mathbf{s}^{(t)}\|_1 \right] + \mu \|\mathbf{L}\|_F^2,$$

---

**Algorithm 1** PG-ELLA ( $k, \lambda, \mu$ )

---

```
1:  $T \leftarrow 0$ 
2:  $\mathbf{A} \leftarrow \text{zeros}_{k \times d, k \times d}$ ,  $\mathbf{b} \leftarrow \text{zeros}_{k \times d, 1}$ 
3:  $\mathbf{L} \leftarrow \text{RandomMatrix}_{d, k}$ 
4: while some task  $(\mathcal{Z}^{(t)}, \phi(\mathbf{m}^{(t)}))$  is available do
5:   if isNewTask( $\mathcal{Z}^{(t)}$ ) then
6:      $T \leftarrow T + 1$ 
7:      $(\mathbb{T}^{(t)}, R^{(t)}) \leftarrow \text{getRandomTrajectories}()$ 
8:   else
9:      $(\mathbb{T}^{(t)}, R^{(t)}) \leftarrow \text{getTrajectories}(\alpha^{(t)})$ 
10:     $\mathbf{A} \leftarrow \mathbf{A} - (\mathbf{s}^{(t)} \mathbf{s}^{(t)\top}) \otimes \mathbf{\Gamma}^{(t)}$ 
11:     $\mathbf{b} \leftarrow \mathbf{b} - \text{vec}(\mathbf{s}^{(t)\top} \otimes (\alpha^{(t)\top} \mathbf{\Gamma}^{(t)}))$ 
12:   end if
13:   Compute  $\alpha^{(t)}$  and  $\mathbf{\Gamma}^{(t)}$  from  $\mathbb{T}^{(t)}$ 
14:    $\mathbf{s}^{(t)} \leftarrow \arg \min_{\mathbf{s}} \left\| \alpha^{(t)} - \mathbf{L} \mathbf{s}^{(t)} \right\|_{\mathbf{\Gamma}^{(t)}}^2 + \mu \|\mathbf{s}\|_1$ 
15:    $\mathbf{A} \leftarrow \mathbf{A} + (\mathbf{s}^{(t)} \mathbf{s}^{(t)\top}) \otimes \mathbf{\Gamma}^{(t)}$ 
16:    $\mathbf{b} \leftarrow \mathbf{b} + \text{vec}(\mathbf{s}^{(t)\top} \otimes (\alpha^{(t)\top} \mathbf{\Gamma}^{(t)}))$ 
17:    $\mathbf{L} \leftarrow \text{mat} \left( \left( \frac{1}{T} \mathbf{A} + \lambda \mathbf{I}_{k \times d, k \times d} \right)^{-1} \frac{1}{T} \mathbf{b} \right)$ 
18:   for  $t \in \{1, \dots, T\}$  do:  $\theta^{(t)} \leftarrow \mathbf{L} \mathbf{s}^{(t)}$ 
19: end while
```

---

where  $\mathbf{S}$  is the set of the sparse vectors  $\mathbf{s}^{(t)}$ , the L1 norm of  $\|\mathbf{s}^{(t)}\|_1$  provides sparse code for  $\mathbf{s}^{(t)}$  and the Frobenious norm in  $\|\mathbf{L}\|_F^2$  provides regularization of  $\mathbf{L}$ . The coefficients  $\mu$  and  $\lambda \in \mathbb{R}$  are weights for the regularization and sparsity respectively. The learning objective function is approximated by a second order Taylor expansion around an estimate  $\alpha^{(t)}$  of the single task policy parameters of task  $t$ . The optimization problem is solved by using the on-line ELLA algorithm introduced in [19] and extended to reinforcement learning in [5]. The optimization problem is solved by incrementally updating  $\mathbf{L}$  by the following the update equations,

$$\begin{aligned} \mathbf{s}^{(t)} &\leftarrow \arg \min_{\mathbf{s}} \left\| \alpha^{(t)} - \mathbf{L} \mathbf{s}^{(t)} \right\|_{\mathbf{\Gamma}^{(t)}}^2 + \mu \|\mathbf{s}\|_1, \\ \mathbf{A} &\leftarrow \mathbf{A} + (\mathbf{s}^{(t)} \mathbf{s}^{(t)\top}) \otimes \mathbf{\Gamma}^{(t)}, \\ \mathbf{b} &\leftarrow \mathbf{b} + \text{vec}(\mathbf{s}^{(t)\top} \otimes (\alpha^{(t)\top} \mathbf{\Gamma}^{(t)})), \\ \mathbf{L} &\leftarrow \text{mat} \left( \left( \frac{1}{T} \mathbf{A} + \lambda \mathbf{I}_{l \times d_{\theta}, l \times d_{\theta}} \right)^{-1} \frac{1}{T} \mathbf{b} \right). \end{aligned} \quad (2)$$

where  $\|\mathbf{v}\|_{\mathbf{A}}^2 = \mathbf{v}^\top \mathbf{A} \mathbf{v}$  and  $\mathbf{\Gamma}^{(t)}$  is the Hessian of the PG objective function, and  $\mathbf{A}$  and  $\mathbf{b}$  are initialized to zero matrices. The full algorithm is described in Algorithm 1.

## 4. EXPERIMENTS

We evaluate our approach by modeling the control policies for different Turtlebot systems [3]. Experiments were conducted using the Gazebo simulator [1, 2] with a simulated model of the Turtlebot as shown in Fig 1. The implementation of our approach uses Python and the Hydro version of Robot Operating System (ROS). This experimental setup allows us to use the same code in both simulation and the physical robots.

In a first attempt to explore disturbance rejection as a preamble of fault tolerant control applications, the Turtlebot should learn how to drive itself from an initial point to a goal

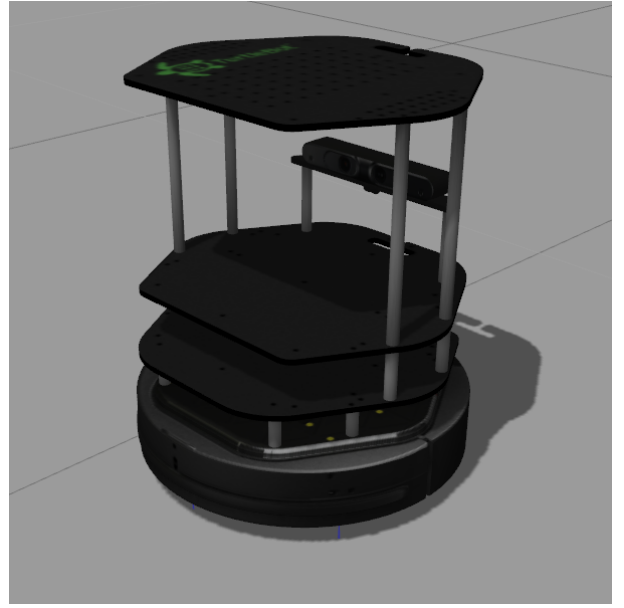


Figure 1: Turtlebot model in Gazebo Simulator.

point. We then artificially induce a random and constant disturbance to the control signal that is drawn uniformly from  $[-0.1, 0.1] \subset \mathbb{R}$  and measured in  $m/s$ . These limits were selected to provide a large noise that was still within the bounds of what the turtlebot accepts. The disturbance now emulates a bias on the angular velocity of each robot that forces the robot to compensate for the induced failure. It is worth mentioning that the difficulty of the disturbance can be increased by introducing time varying disturbance and stochastic disturbances later on.

We assume we have little knowledge of the Turtlebot model, so we use the simplified kinematic model provided in [4]. In the kinematic model, the state space is given by  $\mathbf{x} \in \mathbb{R}^3$  and the action space is described by  $\mathbf{a} \in \mathbb{R}^2$ . Notice that the model in [4] just considers the kinematics of a unicycle in polar coordinates so we do not consider the dynamics of the system, *i.e.*, we are neglecting model parameters such as mass, damping and friction coefficients, as well as inputs such as forces and torques. Then, the nonlinear policy is derived neglecting the dynamics and just assuming simple kinematics, therefore, taking into account that our action is given by  $\mathbf{a} = \theta^\top \phi(\mathbf{x}) = (u, w)^\top$  where  $\mathbf{x} = (\rho, \gamma, \psi)^\top$  as illustrated in Fig. 3 we propose the following nonlinear gain vector structure,

$$\phi(\mathbf{x}) = \begin{pmatrix} \rho \cos(\gamma) \\ \frac{\cos(\gamma) \sin(\gamma)}{\gamma} (\gamma + \psi) \\ 1 \end{pmatrix} \quad (3)$$

where  $\rho$ ,  $\gamma$  and  $\psi$  are indicated in Fig. 2.

The state is a non-linear transformation of the position and heading angle. We ran experiments using natural actor critic [18], episodic REINFORCE [23], and FD [12] and found FD to work best for our problem. We believe [Gabe: Is this speculating?] this is a result of the particular nonlinear transformation we are using and the fact that FD optimizes the true return. In future work we will present a

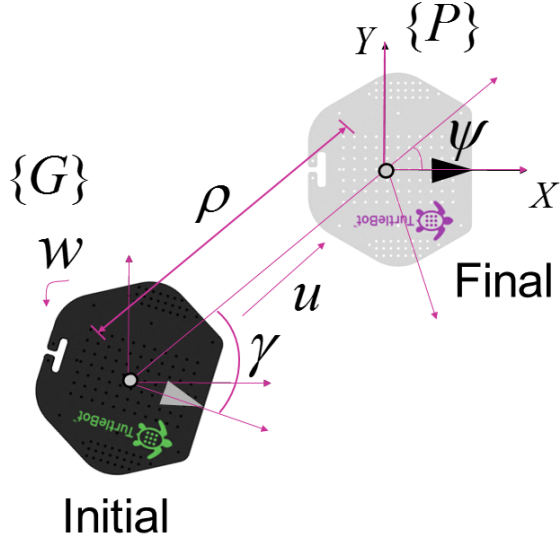


Figure 2: State variables of simplified go-to-goal problem.

comparison of different base learners.

#### 4.1 Methodology

We start by generating 20 robots, each with a different constant disturbance and a unique goal, both selected uniformly. We use 20 robots because it provides a large diversity in tasks but is a number that is still practical to run in simulation.

FD is used to train each robot for 20 iterations with 15 roll-outs per iteration and 50 time-steps per roll-out. Note that 20 iterations is fewer than is required for any system to converge to a good controller. The number of roll-outs and time steps were selected to allow for successful learning while minimizing the runtime. The policy that is learned after 20 iterations is used as  $\alpha^{(t)}$  for PG-ELLA.

We learn our PG-ELLA knowledge repository  $\mathcal{L}$  and sparse representation  $\mathbf{s}^{(t)}$  using the update equations given by (2). Tasks are encountered randomly with repetition and learning stops once every task has been observed once. For our experiments we approximate the hessian with the identity matrix. For the parameters unique to PG-ELLA, we select  $k = 8$  to be the number of columns, and use sparsity coefficient  $\mu = 1 \times 10^{-3}$  and regularization coefficient  $\lambda = 1 \times 10^{-8}$ . These coefficients were selected by trial and error. Other parameters include the learning rate  $\delta = 1 \times 10^{-6}$  and policy noise  $\sigma = 0.001$ .

We then compare the policy reconstructed from PG-ELLA against the policy that was learned after 20 iterations of FD by comparing the learning curves that result from running FD for an additional 80 learning iterations. Performance is averaged over 6 trials for all 20 robots to increase our confidence in the results. In Figure 3 we see that PG-ELLA is successfully able to reconstruct the policies and provide an additional benefit of positive transfer. We plot the change in reward of using PG-ELLA instead of PG in Figure 4, and observe that PG-ELLA improves the quality of the learned controllers. Note that these are preliminary results and we

suspect further refinements will enable us to achieve larger amounts of transfer.

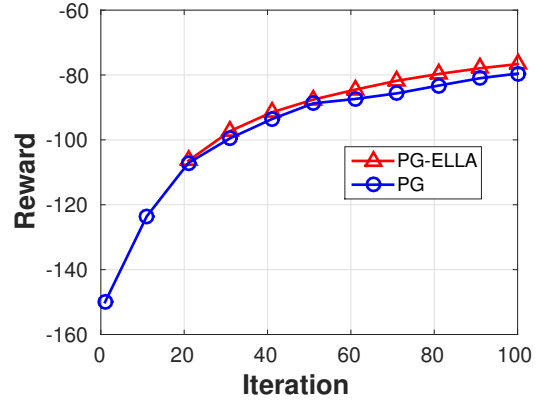


Figure 3: Learning curves for PG and PG-ELLA. Using PG-ELLA to transfer information between tasks improves the performance over standard PG.

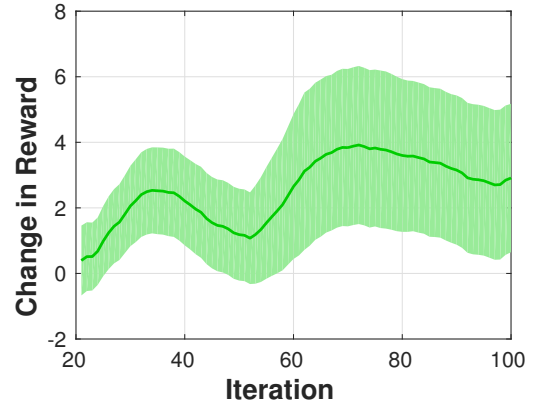


Figure 4: The positive transfer achieved by using ELLA.

## 5. CONCLUSIONS

We demonstrate the use of lifelong learning for disturbance rejection on Turtlebots. This is intended to lay the foundation for creating fault tolerant control on multi-agent systems. We show that PG-ELLA can be implemented on simulated turtlebots and that it outperforms standard PG. This suggests that PG-ELLA can be extended to real robotic systems.

## Acknowledgments

the grant goes here

## 6. ADDITIONAL AUTHORS

## REFERENCES

- [1] Gazebo. [Online]. Available: <http://gazebosim.org/>, 2016.

- [2] Ros.org|powering the world's robots. [Online]. Available: <http://www.ros.org/>, 2016.
- [3] Turtlebot 2. [Online]. Available: <http://www.turtlebot.com/>, 2016.
- [4] M. Aicardi, G. Casalino, A. Balestrino, and A. Bicchi. Closed loop smooth steering of unicycle-like vehicles. In *Proc. of IEEE Conference on Decision and Control*, pages 2455–2458, 1994.
- [5] H. Bou Ammar, E. Eaton, and P. Ruvolo. Online multi-task learning for policy gradient methods. *International Conference on Machine Learning*, 2014.
- [6] H. Bou Ammar, E. Eaton, P. Ruvolo, and M. E. Taylor. Unsupervised cross-domain transfer in policy gradient reinforcement learning via manifold alignment. *International Joint Conference on Artificial Intelligence*, 2015.
- [7] C. G. Cassandras and S. Lafortune. *Introduction to Discrete Event Systems*. Springer, New York, 2nd edition, 2008.
- [8] P. Dorato, C. Abdallah, and V. Cerone. *Linear Quadratic Control*. Krieger, 2000.
- [9] H. Khalil. *Nonlinear Systems*. Prentice Hall, 2002.
- [10] A. Kleiner, M. Dietl, and B. Nebel. Towards a life-long learning soccer agent. In *RoboCup 2002: Robot Soccer World Cup VI*, pages 126–134. Springer, 2002.
- [11] J. Kober, J. A. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, page 0278364913495721, 2013.
- [12] J. Kober, J. A. D. Bagnell, and J. Peters. Reinforcement learning in robotics: A survey. *International Journal of Robotics Research*, July 2013.
- [13] J. Kober and J. Peters. Policy search for motor primitives in robotics. *Advances in Neural Information Processing Systems*, pages 849–856, 2009.
- [14] A. Kumar and H. Daume III. Learning task grouping and overlap in multi-task learning. *arXiv preprint arXiv:1206.6417*, 2012.
- [15] F. L. Lewis and V. L. Syrmos. *Optimal Control*. John Wiley & Sons, 3rd edition, 2012.
- [16] J. M. Luna, C. T. Abdallah, and G. Heileman. Performance optimization and regulation for multitier servers. In *Proc. of IEEE International Conference on Decision and Control*, pages 1026–1032, Osaka, December 2015.
- [17] N. S. Nise. *Control Systems Engineering*. John Wiley & Sons, 7th edition, 2010.
- [18] J. Peters and S. Schaal. Natural actor-critic. *Neurocomputing*, 71(7):1180–1190, 2008.
- [19] P. Ruvolo and E. Eaton. ELLA: An efficient lifelong learning algorithm. *International Conference on Machine Learning*, 28:507–515, 2013.
- [20] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in Neural Information Processing Systems*, 99:1057–1063, 1999.
- [21] S. Thrun and T. M. Mitchell. *Lifelong robot learning*. Springer, 1995.
- [22] B. Ugaonkar, G. Pacifi, P. Shenoy, M. Spreitzer, and A. Tantawi. Analytic modeling of multitier internet applications. *ACM Trans. on the Web*, 1(1), May 2007.
- [23] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.