

Grounded Action Transformation for Robot Learning in Simulation

Josiah P. Hanna and Peter Stone

{jphanna,pstone}@cs.utexas.edu
The University of Texas at Austin
2317 Speedway, Austin, TX 78712 USA

Abstract

Robot learning in simulation is a promising alternative to the prohibitive sample cost of learning in the physical world. Unfortunately, policies learned in simulation often perform worse than hand-coded policies when applied on the physical robot. This paper proposes a new algorithm for learning in simulation – Grounded Action Transformation – and applies it to learning of humanoid bipedal locomotion. Our approach results in a 43.27% improvement in forward walk velocity compared to a state-of-the-art hand-coded walk.¹

Introduction

An alternative to learning robotic skills directly on the physical robot is learning in simulation and transferring the learned skills. However, the value of simulation learning is limited by the inherent inaccuracy of simulators in modeling the dynamics of the physical world. As a result, learning that takes place in a simulator is unlikely to increase real world performance.

Grounded Simulation Learning (GSL) is a framework for modifying a simulator such that learning transfers to the physical robot (Farchy et al. 2013). We introduce a new GSL algorithm – Grounded Action Transformation (GAT) for simulation-transfer. The algorithm is evaluated on the task of bipedal robotic walking. Our evaluation of GAT starts from a state-of-the-art walk engine as the base policy and improves the walk velocity by over 43%, leading to what may be the fastest known walk on the SoftBank NAO robot.

Problem Statement

At discrete time-step t the robot takes action $A_t \sim \pi(\cdot|S_t)$ according to π which is a distribution over actions, $A_t \in \mathcal{A}$, conditioned on the current state, $S_t \in \mathcal{S}$. The environment, E , responds with $S_{t+1} \sim P(\cdot|S_t, A_t)$ according to the dynamics, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}_{\geq 0}$ which is a probability density function over next states conditioned on the current state and action. A trajectory of length L is a sequence of states and

actions, $\tau := S_0, A_0, \dots, S_L, A_L$. The cost function, c , assigns a scalar cost to each (s, a) pair. The value of policy, π , is defined as $J(\pi) := \mathbb{E}_{\tau \sim Pr(\cdot|\pi)} \left[\sum_{t=0}^L c(S_t, A_t) \right]$ where $Pr(\tau|\pi)$ is the probability of τ when selecting actions according to π . The policy π is parameterized by a vector θ . Since θ determines the policy distribution we overload notation and refer to π_θ as θ . Given initial parameters θ_0 , the goal of policy improvement is to find θ' such that $J(\theta') < J(\theta_0)$.

In this paper, learning takes place in a simulator, E_{sim} , that models E . Specifically E_{sim} has the same state-action space as E but inevitably a different dynamics distribution, P_{sim} . The different dynamics distribution mean that for any policy, θ , $J(\theta) \neq J_{sim}(\theta)$ since θ induces a different trajectory distribution in E than in E_{sim} . Thus θ' with $J_{sim}(\theta') < J_{sim}(\theta_0)$ does not imply $J(\theta') < J(\theta_0)$ – in fact $J(\theta')$ could even be worse than $J(\theta_0)$. This paper explores methods for learning in E_{sim} that lower expected cost.

Grounded Action Transformation

The Grounded Simulation Learning (GSL) framework (Farchy et al. 2013) improves robot learning in simulation by using state transition data from the physical system to ground E_{sim} such that it is a better model of E . The GSL framework assumes there is an imperfect simulator, $E_{sim} = \langle \mathcal{S}, \mathcal{A}, P_{sim}, c \rangle$, that models E . Furthermore, E_{sim} must be *modifiable*. A modifiable simulator has parameterized transition probabilities $P_\phi(\cdot|s, a) := P_{sim}(\cdot|s, a; \phi)$ where the vector ϕ can be changed to produce, in effect, a different simulator.

GSL grounds E_{sim} by finding ϕ^* such that:

$$\phi^* = \operatorname{argmin}_{\phi} \sum_{\tau \in \mathcal{D}} d(Pr(\tau|\theta), Pr_{sim}(\tau|\theta, \phi)) \quad (1)$$

where $Pr(\tau|\theta)$ is the probability of observing trajectory τ in E , $Pr_{sim}(\tau|\theta, \phi)$ is the probability of τ in E_{sim} with dynamics parameterized by ϕ , d is a measure of similarity between these probabilities, and \mathcal{D} is a data set of trajectories recorded on the physical robot. When a policy improvement method is used to optimize θ in E_{sim} with ϕ^* the improved policy, θ' is expected to perform better on the physical robot.

Our principle algorithmic contribution, GSL with *Grounded Action Transformation* (GAT) improves grounded by correcting differences in the effects of actions between

Copyright © 2017, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹An extended version of this work is published in AAAI 2017. Videos of learned walk policies can be found at https://www.cs.utexas.edu/users/AustinVilla/?p=research/real_and_sim_walk_learning.

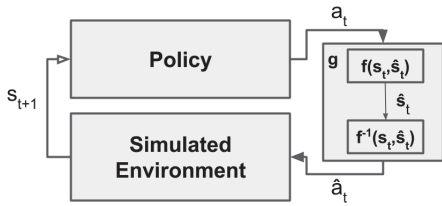


Figure 1: GAT induces a modifiable simulator.

E and E_{sim} . Since it is often easier to minimize error in the one step dynamics distribution, GAT uses:

$$\phi^* = \underset{\phi}{\operatorname{argmin}} \sum_{\tau_i \in \mathcal{D}} \sum_{t=0}^L d(P(s_{t+1}^i | s_t^i, \mathbf{a}_t^i), P_{\phi}(s_{t+1}^i | s_t^i, \mathbf{a}_t^i)) \quad (2)$$

as a surrogate loss function for (1) which can be minimized with transitions observed in \mathcal{D} .

To find ϕ^* efficiently, GAT uses a parameterized *action transformation* function which takes the agent’s state and action as input and outputs a new action which – when taken in simulation – will result in the robot transitioning to the same next state it would have in E . We denote this function, $g_{\phi} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{A}$; the parameters of g serve as ϕ under the GSL framework. GAT constructs g with parameterized models of the robot’s dynamics and the simulators *inverse* dynamics. Assuming a dataset \mathcal{D}_{sim} of simulated experience, GAT reduces (2) to a supervised learning problem.

GAT defines $g := g_{\phi}$ by a deterministic forward model of the robot’s dynamics, f and a deterministic model of the simulator’s inverse dynamics, f_{sim}^{-1} . Let \mathbf{x}_t be the robot’s joint configuration at time t . The function, f maps (s_t, \mathbf{a}_t) to the maximum likelihood estimate of \mathbf{x}_{t+1} under P . The function f_{sim}^{-1} maps (s_t, s_{t+1}) to the action that is most likely to produce this transition in simulation. When executing θ in simulation, the robot selects $\mathbf{a}_t \sim \pi_{\theta}(\cdot | s_t)$ and then uses f to predict what the resulting configuration, \mathbf{x}_{t+1} , would be in E . Then \mathbf{a}_t is replaced with $\hat{\mathbf{a}}_t := f_{\text{sim}}^{-1}(s_t, f(s_t, \mathbf{a}_t))$. The result is the robot achieves the exact \mathbf{x}_{t+1} it would have on the physical robot. In practice f and f_{sim}^{-1} are represented with supervised regression models (neural networks in this work) and learned from \mathcal{D} and \mathcal{D}_{sim} respectively. Figure 1 illustrates the GAT-modified E_{sim} . GAT then proceeds to improve θ with optimize within the grounded simulator.

Empirical Results

We evaluate GAT on the task of bipedal robot walking using the SoftBank NAO. For walking, our NAO uses an open source walk engine developed at the University of New South Wales (UNSW) (Ashar et al. 2015). This walk engine has been used by at least one team in each of the past three RoboCup Standard Platform League (SPL) championship games in which teams of five NAOs compete in soccer matches. On the physical robot a trajectory terminates once the robot has walked four meters or falls. A trajectory generated with θ_0 lasts ≈ 20.5 seconds on the robot. We use two different simulators in this work: the open source SimSpark and Gazebo

θ_0	GAT with SimSpark	GAT with Gazebo
19.52 (cm/s)	26.27	26.89
0.0 (%)	34.58	37.76

Table 1: This table gives the maximum learned velocity and percent improvement over θ_0 (left column).

simulators. In simulation a trajectory terminates after a fixed time. We use the Covariance Matrix Adaption-Evolutionary Strategy (CMA-ES) algorithm (Hansen 2006) for policy improvement within simulation. To clarify terminology, a generation refers to a single update of CMA-ES; an iteration refers to a complete cycle of GSL.

We evaluate GAT for transferring simulation learning to the physical NAO using both SimSpark and Gazebo as E_{sim} . The data set \mathcal{D} consists of 15 trajectories collected with θ_0 on the physical NAO. For each iteration, we optimize θ for 10 generations of the CMA-ES algorithm. We evaluate the best policy from each generation with 5 trajectories on the physical robot. If the robot falls in any of the 5 trajectories the policy is considered unstable.

Experimental Results

Table 1 gives the walk velocity of the best θ' returned by GAT in each simulator. GAT with SimSpark and GAT with Gazebo both improved walk velocity by over 30%. This result demonstrates generality of GAT across different simulators. Two iterations of GAT with SimSpark (not shown) increased the walk velocity of the NAO by 43.27% compared to θ_0 .

Policy improvement with CMA-ES required 30,000 trajectories per iteration to find the 10 policies that were evaluated on the robot. In contrast the total number of trajectories executed on the physical robot is 65 (15 trajectories in \mathcal{D} and 5 evaluations per θ'). This result demonstrates GAT can use sample-intensive simulation learning to optimize real world skills with a low number of trajectories on the physical robot.

Conclusion

We proposed and evaluated the Grounded Action Transformation (GAT) algorithm for grounded simulation learning. Our method led to a 43.27 % improvement in the walk velocity of a state-of-the-art bipedal robot. Furthermore, this experiment demonstrated the benefits of GAT are independent of the choice of simulator.

References

- Ashar, J.; Ashmore, J.; Hall, B.; and Harris, S. e. a. 2015. Robocup spl 2014 champion team paper. In *RoboCup 2014: Robot World Cup XVIII*, volume 8992 of *Lecture Notes in Computer Science*. Springer International Publishing. 70–81.
- Farchy, A.; Barrett, S.; MacAlpine, P.; and Stone, P. 2013. Humanoid robots learning to walk faster: From the real world to simulation and back. In *Twelfth International Conference on Autonomous Agents and Multiagent Systems, AAMAS*.
- Hansen, N. 2006. The CMA evolution strategy: a comparing review. In Lozano, J.; Larranaga, P.; Inza, I.; and Bengoetxea, E., eds., *Towards a new evolutionary computation. Advances on estimation of distribution algorithms*. Springer. 75–102.