

Bipedal walking using deep reinforcement learning

Jhon Charaja¹ and Luca Borgonovi¹

¹Universidade de São Paulo, Brasil

December 10, 2021

Outline

1. Motivation

2. Objective

3. Scopes

4. Methodology

5. Results

6. Conclusions

Motivation

Bipedal walking is difficult activity to describe due to nonlinear relationships.

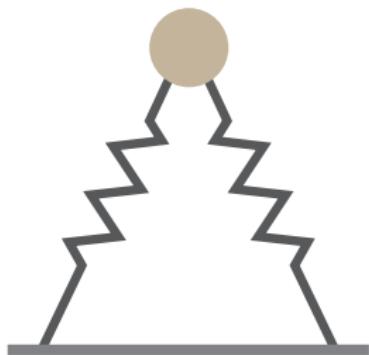


Figure: Spring-loaded inverted pendulum

- Physical representation of legs as springs¹
- Motion of the trunk with center of mass¹
- Recent works perform bipedal walking with machine learning methods^{2,3}
- Deep reinforcement learning (DRL) does not require mathematical model or control formulation⁴
- DRL with proximal policy optimization is simple and efficient⁴

¹H. Geyer (2006).

²T. Haarjona (2019).

³T. Li (2019).

⁴J. Schulman (2017)

Motivation

Bipedal walking is difficult activity to describe due to nonlinear relationships.

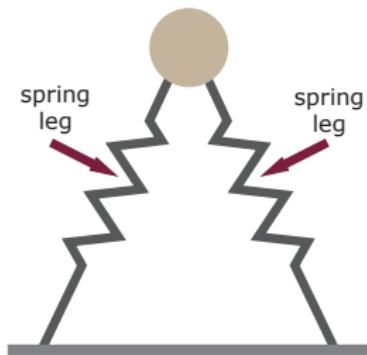


Figure: Spring-loaded inverted pendulum

- Physical representation of legs as springs¹
- Motion of the trunk with center of mass¹
- Recent works perform bipedal walking with machine learning methods^{2,3}
- Deep reinforcement learning (DRL) does not require mathematical model or control formulation⁴
- DRL with proximal policy optimization is simple and efficient⁴

¹H. Geyer (2006).

²T. Haarjona (2019).

³T. Li (2019).

⁴J. Schulman (2017)

Motivation

Bipedal walking is difficult activity to describe due to nonlinear relationships.

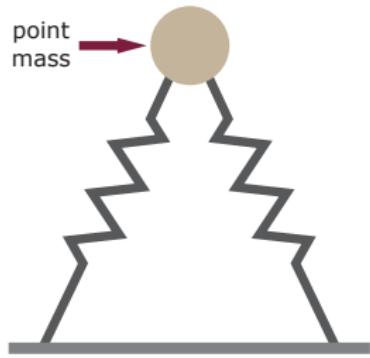


Figure: Spring-loaded inverted pendulum

- Physical representation of legs as springs¹
- Motion of the trunk with center of mass¹
- Recent works perform bipedal walking with machine learning methods^{2,3}
- Deep reinforcement learning (DRL) does not require mathematical model or control formulation⁴
- DRL with proximal policy optimization is simple and efficient⁴

¹H. Geyer (2006).

²T. Haarjona (2019).

³T. Li (2019).

⁴J. Schulman (2017)

Motivation

Bipedal walking is difficult activity to describe due to nonlinear relationships.

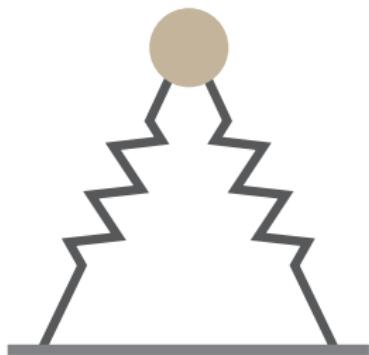


Figure: Spring-loaded inverted pendulum

- Physical representation of legs as springs¹
- Motion of the trunk with center of mass¹
- Recent works perform bipedal walking with machine learning methods^{2,3}
- Deep reinforcement learning (DRL) does not require mathematical model or control formulation⁴
- DRL with proximal policy optimization is simple and efficient⁴

¹H. Geyer (2006).

²T. Haarjona (2019).

³T. Li (2019).

⁴J. Schulman (2017)

Outline

1. Motivation

2. Objective

3. Scopes

4. Methodology

5. Results

6. Conclusions

Perform bipedal walking using deep reinforcement learning with proximal policy optimization algorithm

Outline

1. Motivation

2. Objective

3. Scopes

4. Methodology

5. Results

6. Conclusions

Scopes

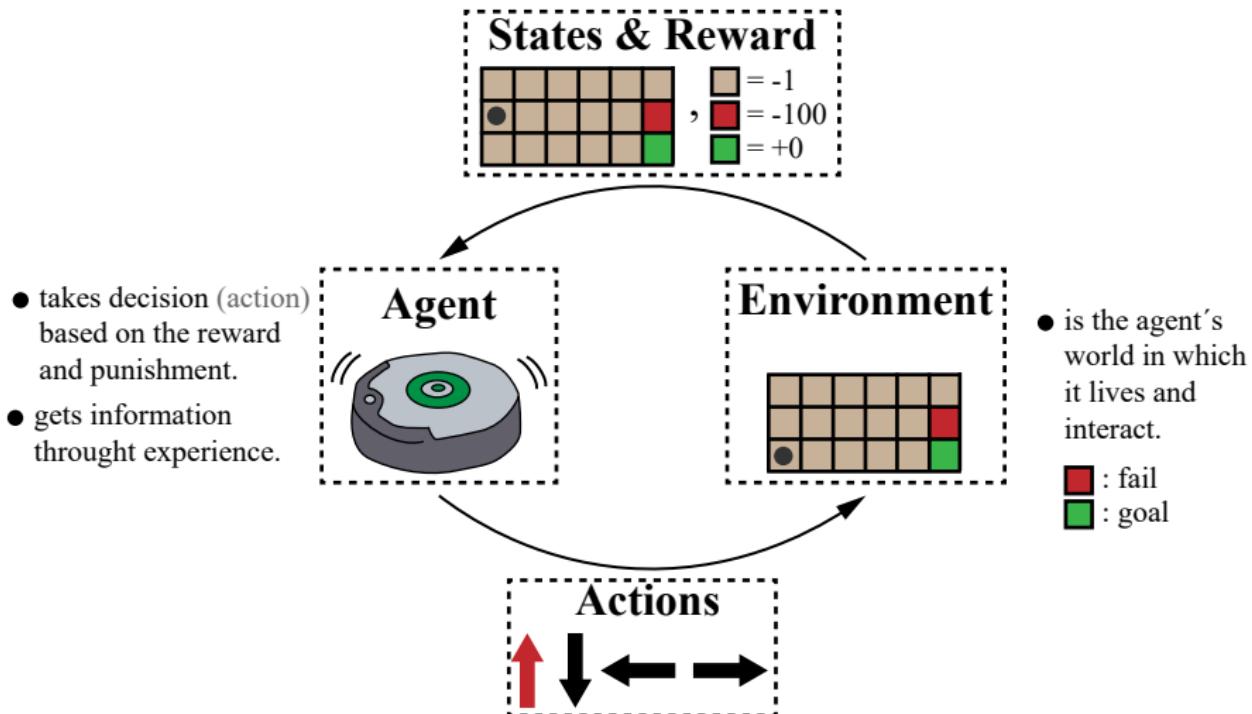
- The algorithm is focused on making the robot walk, but not on controlling the way it walks.
- Likewise, the algorithm is not focused on optimizing energy during walking.

Outline

1. Motivation
2. Objective
3. Scopes
4. Methodology
5. Results
6. Conclusions

What is reinforcement learning?

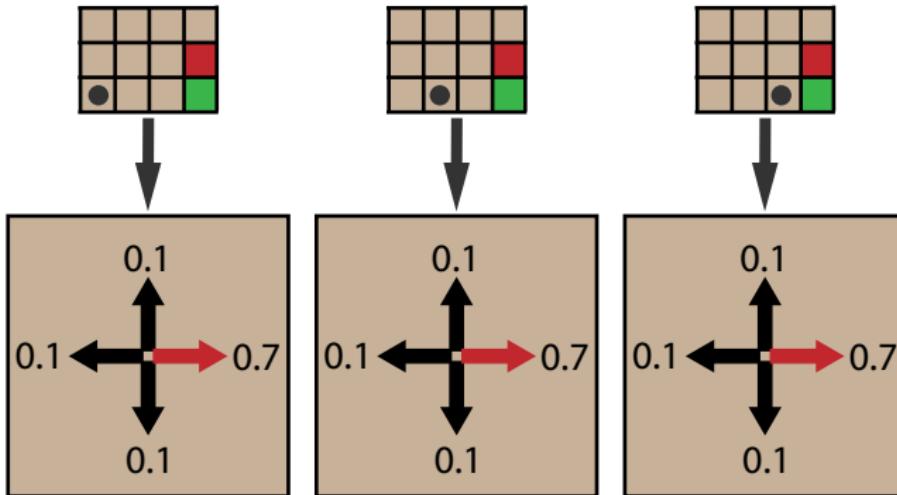
Reinforcement learning is a training technique of machine learning to teach models (agents) to make good (optimal) sequences of decisions under uncertainty¹.



¹R. S. Sutton (2018).

Policy

A policy (π) indicates the decision (action) that the agent should take as a function of the agent's state to accomplish the task¹.



The objective is find the optimal policy (π) that maximizes the reward

How quantify the performance of a policy (π)?

¹R. S. Sutton (2018).

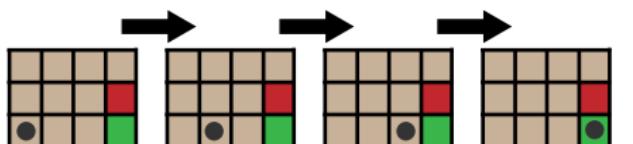
Policy evaluation: value function

The value function ($V^\pi(s)$) indicates the final return from being in a state when following a particular policy (π)¹. It is given by

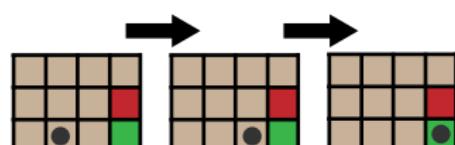
$$V^\pi(s_t = s) = E_\pi [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s],$$

where γ is a discount factor and s_t is initial state.

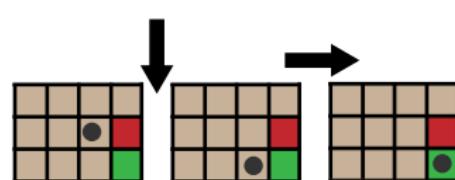
Policy (π): Considering: $\gamma=1$



$$V^\pi(s_1): -2$$



$$V^\pi(s_2): -1$$



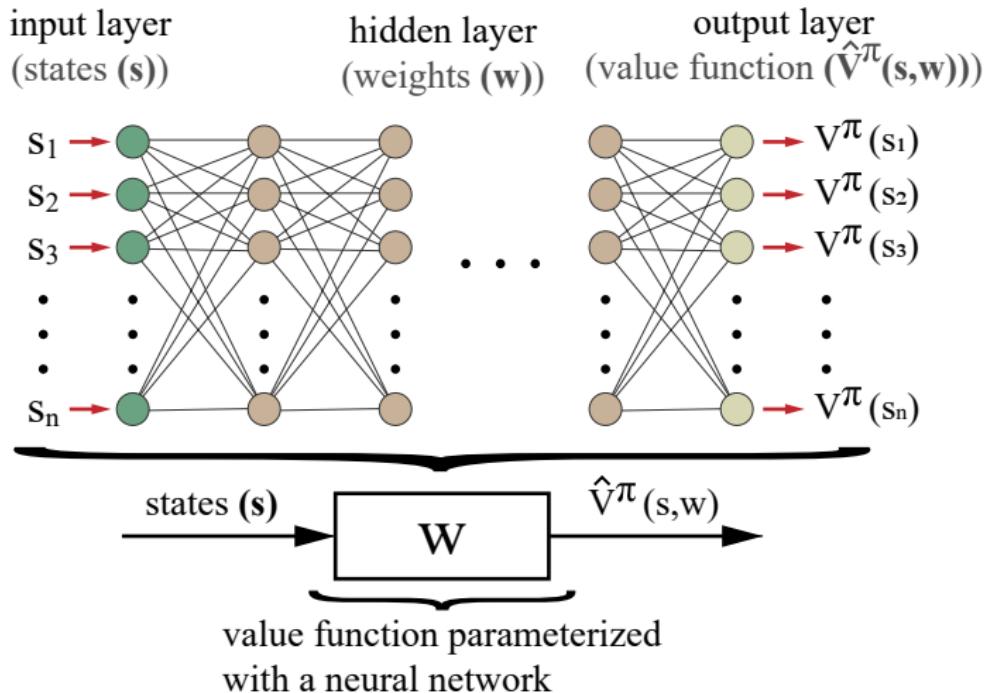
$$V^\pi(s_6): -1$$

- $V^\pi(s)$ depends of policy (π) and initial state (s_t).
- Agent needs to start in $(s_{3,4,5,7,8,9,10})$ to know its value $V^\pi(s)$.

¹R. S. Sutton (2018).

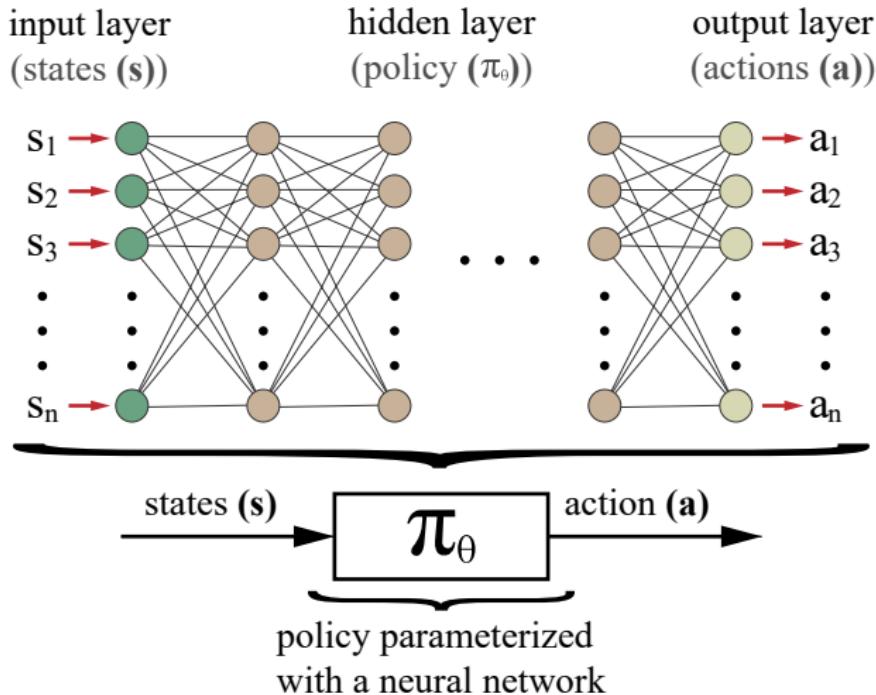
Value function approximation

The value function could be approximated by a neural network



Policy and neural networks

The policy could be represented by a neural network



Policy gradient method

This method is focused on modify the neural network parameters (θ) to get a optimal (local) policy¹.

The objective function (reward) is formulated as¹

$$L(\theta) = \hat{\mathbb{E}}_t \left[\log \pi_\theta(a|s) \hat{A}_t \right],$$

with,

$$\hat{A}_t = V_t^{\pi_\theta} - \hat{V}_t,$$

where,

- \hat{A}_t advantage function
- $V_t^{\pi_\theta}$ final reward following policy (π_θ)
- \hat{V}_t expected final reward (baseline)

¹R. S. Sutton (2018).

Proximal policy optimization

PPO is one of the best algorithms for reinforcement learning due to simplicity and high performance¹.

The main objective function (reward) is given by

$$L^{\text{CLIP}}(\theta) = \hat{\mathbb{E}}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right],$$

with,

$$r_t(\theta) = \frac{\pi_\theta(a|s)}{\pi_{\theta,\text{old}}(a|s)},$$

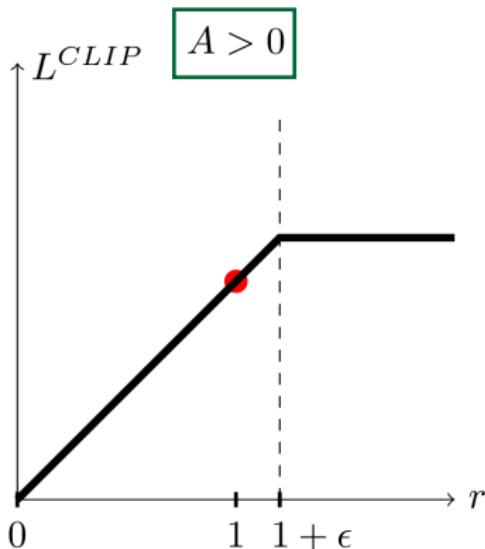
where,

- $r_t(\theta)$ probability ratio between policies
- ϵ hyperparameter

¹J. Schulman (2017).

Proximal policy optimization

the new action yielded
better than expected return



the new action yielded
worse than expected return

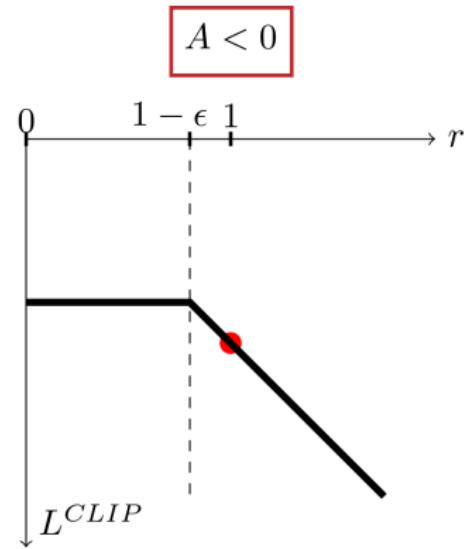


Figure: Image adapted from J. Schulman (2017)

Proximal policy optimization

The final objective function (reward and exploration) is given by¹

$$L^{\text{PPO}}(\theta) = \hat{\mathbb{E}}_t [L^{\text{CLIP}}(\theta) - c_1 L^{\text{VF}}(\theta) + c_2 S[\pi_\theta](s_t)] ,$$

with,

$$L^{\text{VF}}(\theta) = (V_\theta(s_t) - V_t^{\text{target}})^2 ,$$

where,

- c_1, c_2 weighting coefficients
- S entropy bonus (exploration)
- $L^{\text{VF}}(\theta)$ square-error loss
- V_t^{target} objective final reward

¹J. Schulman (2017).

MuJoCo

Advanced physics simulation

OpenAI



TensorFlow

Bipedal robot: algorithm configuration

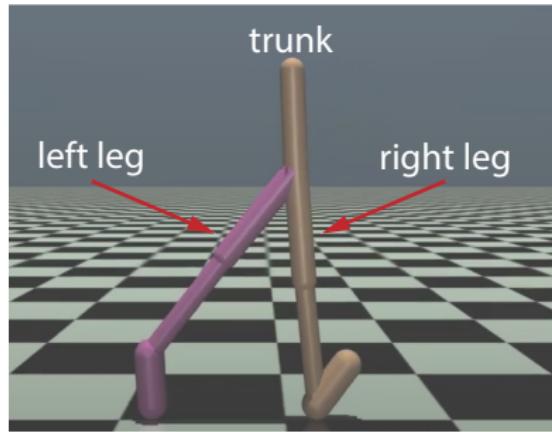


Figure: Bipedal robot in simulation environment

States

- position: 6 (left, right) + 1 (trunk)
- velocity: 6 (left, right) + 1 (trunk)
- total: 14 states

Reward

- fail when trunk angle > 50°
- +1 for each iteration alive
- + linear velocity

Outline

1. Motivation

2. Objective

3. Scopes

4. Methodology

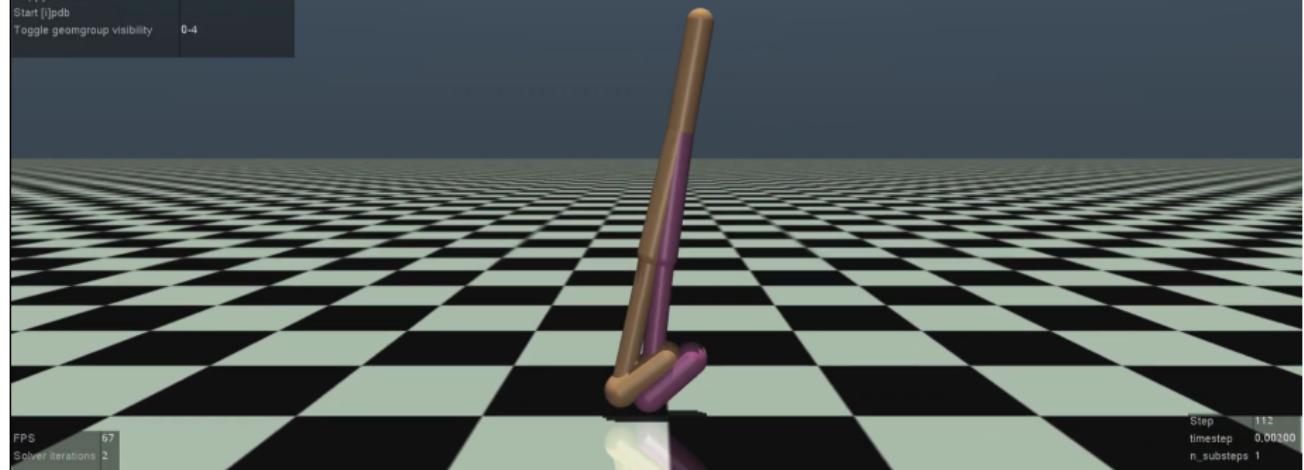
5. Results

6. Conclusions

Results

Run speed = 0.125 x real time
[S]lower, [F]aster
Ren[der] every frame
On
Switch camera (#cams = 2)
[Tab] (camera ID = 0)
[C]ontact forces
On
Referenc[e] frames
On
Tr[ansparent]
Off
Display [M]jocap bodies
On
Stop
[Space]
Advance simulation by one step [right arrow]
[H]ide Menu
Record [V]ideo (Off)
Cap[ture] frame
Start [J]pdb
Toggle geomgroup visibility 0-4

State: 20M
10 attempts per state



Outline

1. Motivation

2. Objective

3. Scopes

4. Methodology

5. Results

6. Conclusions

Conclusions