



## Balancing Safety and Exploration in Policy Gradient

Supervisor: Prof. Marcello Restelli

Co-supervisor: Dott. Matteo Papini

Andrea Battistello, 873795

July 25, 2018

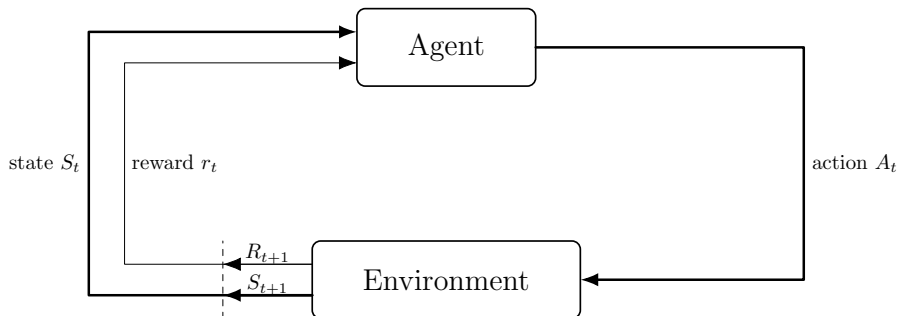
# Contents

- 1 Reinforcement Learning
- 2 Safe Reinforcement Learning
- 3 The role of exploration
- 4 Safely-Exploring Policy Gradient (SEPG)
- 5 Results
- 6 Conclusions

# Reinforcement Learning

We learn from the **interaction** with the environment

For each step, the agent performs an action and receives an observation



# Example: Mountain Car

# A Reinforcement Learning method

## Policy

A policy  $\pi_{\theta}$  is a function that maps states to actions

## Performance

The performance  $J(\theta)$  of a policy is the discounted sum of rewards obtained by following  $\pi_{\theta}$

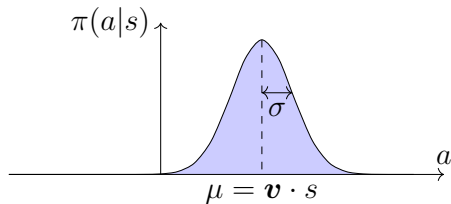
$$J(\theta) = \mathbb{E}_{a_t \sim \pi_{\theta}} \left[ \sum_{k=0}^T \gamma^k r(s_k, a_k) \right]$$

## Policy Gradient method

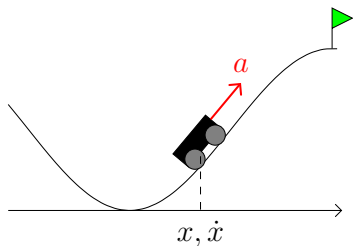
The optimal policy can be found by gradient ascent on policy parameters:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

# Gaussian policies



$$\pi(a|s) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \left(\frac{a - \mathbf{v} \cdot \mathbf{s}}{\sigma}\right)^2\right)$$

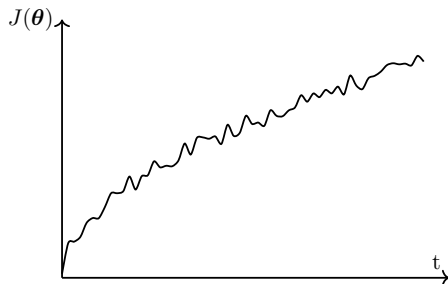


$$\mathbf{s} = [x; \dot{x}]^T$$

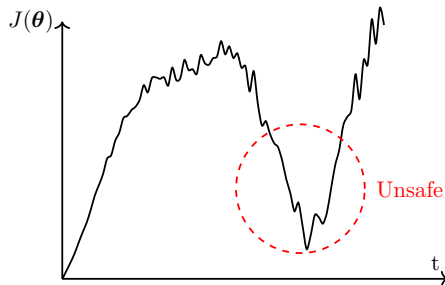
$$\sigma = e^w$$

$$\boldsymbol{\theta} = [\mathbf{v}^T; w]^T$$

# Safe Reinforcement Learning



(a) Ideal behavior



(b) More realistic scenario

# Safe Reinforcement Learning

## Definition

Given a reference  $\underline{J}$  and a confidence level  $\delta$ , a RL algorithm is **safe** if it yields a policy with performance lower than  $\underline{J}$  with probability at most  $\delta$ .

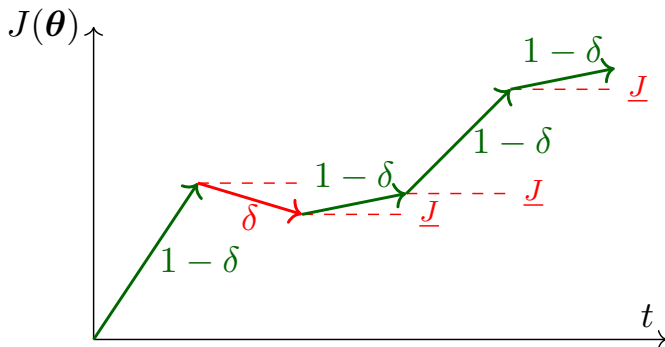


Figure: Example of monotonic improvements



# State of the art in Safe Reinforcement Learning

We will mainly refer on SPG algorithm from (Adaptive Step Size in Policy Gradient, Pirotta et al., 2013).

- Assuming a policy gradient method with Gaussian-parameterized policies  $\pi_{\theta}$ , an update rule of the form  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$  and **fixed variance**  $\sigma = e^w = \text{const}$ :

$$J(\theta') - J(\theta) \geq \alpha \|\nabla_{\theta} J(\theta)\|_2^2 - \alpha^2 \|\nabla_{\theta} J(\theta)\|_1^2 \frac{c_3 + c_2 \sigma}{c_1 \sigma^3}$$

- that is maximized by the following step size:

$$\alpha^* = \frac{c_1 \sigma^3 \|\nabla_{\theta} J(\theta)\|_2^2}{(c_2 \sigma + c_3) \|\nabla_{\theta} J(\theta)\|_1^2},$$

- that guarantees an improvement of:

$$J(\theta') - J(\theta) \geq \frac{1}{2} \alpha^* \|\nabla_{\theta} J(\theta)\|_2^2.$$

# Problems of Safe Reinforcement Learning

The results seen so far suffer from the following problems:

- They target only **monotonic improvements**
  - This covers only specific user needs.
- They are **overly conservative**
  - Resulting in a very slow convergence speed.
- They do not consider an **adaptive variance**
  - The exploration factor remains constant and highly depends on the domain.

# Exploration in Reinforcement Learning

# Contributions (1)

We extended the performance improvement bounds to include an adaptive exploration factor:

$$J(\mathbf{v}, w') - J(\mathbf{v}, w) \geq \beta \nabla_w J(\mathbf{v}', w)^2 - d \beta^2 \nabla_w J(\mathbf{v}', w)^2$$

that is maximized by:

$$\beta^* = \frac{1}{2d}$$

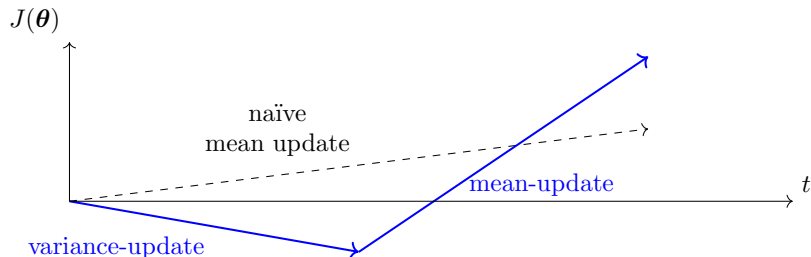
which guarantees:

$$J(\mathbf{v}, w') - J(\mathbf{v}, w) \geq \frac{\nabla_w J(\mathbf{v}', w)^2}{4d}$$

**Recall:**

$$\theta = [\mathbf{v}; w]$$

# Contributions (1)



We adapted the results by employing a new gradient direction:

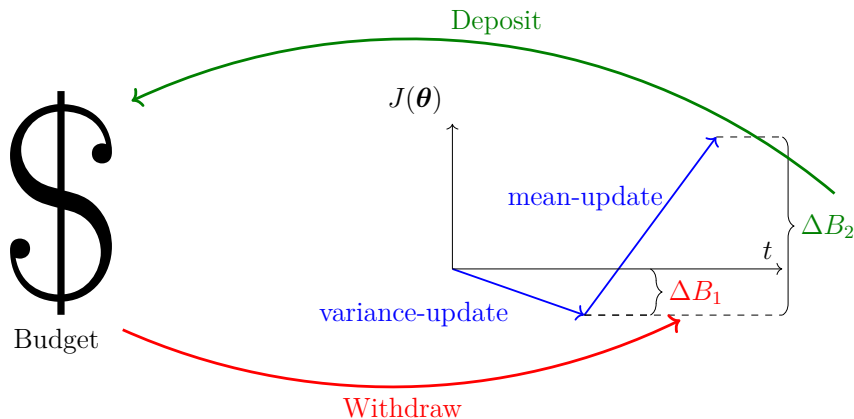
$$w \leftarrow w + \beta \nabla_w \mathcal{L}(\mathbf{v}, w),$$

where:

$$\nabla_w \mathcal{L}(\mathbf{v}, w) := \nabla_w (J(\mathbf{v}', w) - J(\mathbf{v}, w))$$

# Contributions (1)

We employed the budget trick to allow for a customizable implementation of this type of update (and more).



The second contribution was to introduce a new framework that goes beyond the single monotonic improvement case.

## Type-I

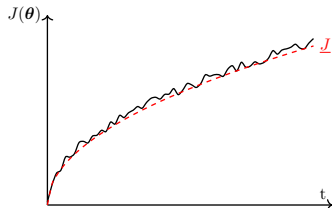
- We guarantee  $J(\theta^t) \geq J_B^t$  for each **policy update**.
- Suited for systems that requires high reliability.

## Type-II

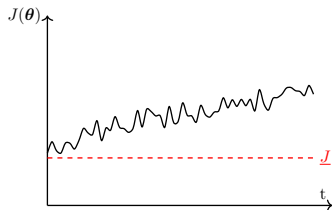
- We guarantee  $J(\theta^t) \geq J_B^t$  **on average** over a learning iteration (e.g., a production day).
- Suited for systems with lower safety needs.

# Contributions (2)

## Type-I

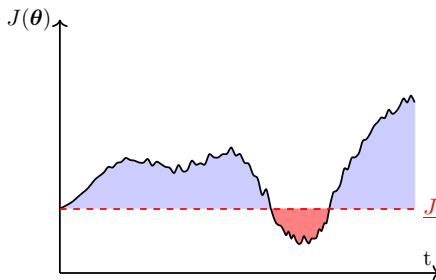


(a) Monotonic improvement (MI)



(b) Lower-bounded I (LB-I)

## Type-II



(c) Lower-bounded II (LB-II)



# Contributions (3)

We devised a new general algorithm that can be customized to specific user needs.

---

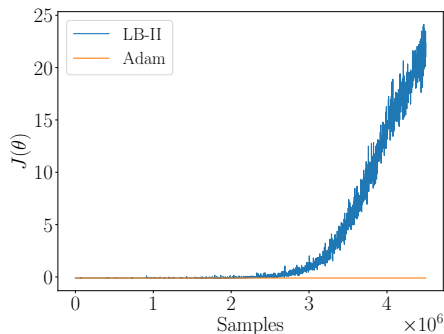
**Algorithm 1** Safely-Exploring Policy Gradient

---

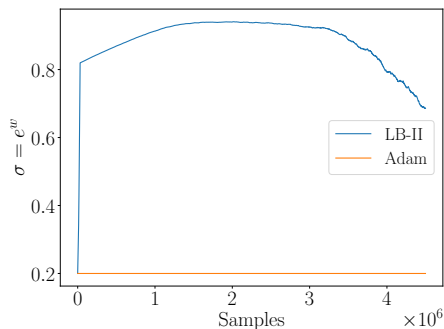
```
1: input:  $\theta^0 = [\mathbf{v}^0, w^0]$ ,  $B^0$ 
2: for  $t = 1, 2 \dots$  do
3:    $\mathbf{v}^{t+1} \leftarrow \mathbf{v}^t + \bar{\alpha} \nabla_{\mathbf{v}} J(\mathbf{v}^t, w^t)$  ▷  $\mathbf{v}$ -update
4:    $B \leftarrow B + J(\mathbf{v}^{t+1}, w^t) - J(\mathbf{v}^t, w^t)$ . ▷  $\mathbf{v}$ -budget update
5:    $w^{t+1} \leftarrow w^t + \bar{\beta} \nabla_w \mathcal{L}(\mathbf{v}^{t+1}, w^t)$  ▷  $w$ -update
6:    $B \leftarrow B + J(\mathbf{v}^{t+1}, w^{t+1}) - J(\mathbf{v}^{t+1}, w^t)$ . ▷  $w$ -budget update
7: end for
```

---

# Results - Mountain Car task

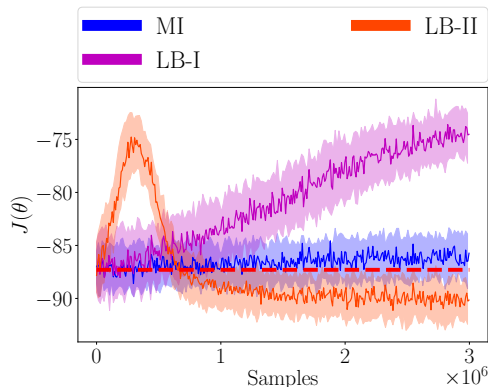


(a) Evolution of the performance

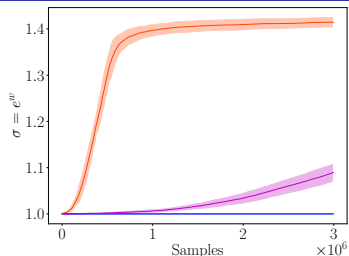


(b) Evolution of the exploration parameter

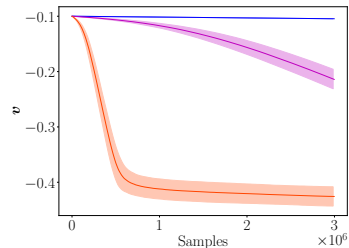
# Results - Linear Quadratic Gaussian controller



(a) Performance of MI, LB-I and LB-II variants



(b) Evolution of variance  $\sigma$



(c) Evolution of the mean  $\nu$

In this work we have introduced:

- An **adaptive way to explore** the environment
- A **new general framework** for safe reinforcement learning
- A **general algorithm** that can be customized to the user needs

Further works can focus on:

- New methods to adapt exploration using previous knowledge about the environment
- New ways to invest the budget
- Extend the result beyond the policy gradient method

Thank you for your attention

Questions?