



**FACULTY
OF MATHEMATICS
AND PHYSICS**
Charles University

BACHELOR THESIS

Michal Pospěch

Thesis title is N/A

Department of Theoretical Computer Science and Mathematical Logic

Supervisor of the bachelor thesis: **Mgr. Roman Neruda, CSc.**

Study programme: **Computer science**

Study branch: **General computer science**

Prague **2021**

I declare that I carried out this bachelor thesis independently, and only with the cited sources, literature and other professional sources. It has not been used to obtain another or the same degree.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Sb., the Copyright Act, as amended, in particular the fact that the Charles University has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 subsection 1 of the Copyright Act.

In date

Author's signature

Dedication. It is nice to say thanks to supervisors, friends, family, book authors and food providers.

Title: Thesis title is N/A

Author: Michal Pospěch

Department: Department of Theoretical Computer Science and Mathematical Logic

Supervisor: Mgr. Roman Neruda, CSc., Department of Theoretical Computer Science and Mathematical Logic

Abstract: Abstracts are an abstract form of art. Use the most precise, shortest sentences that state what problem the thesis addresses, how it is approached, pinpoint the exact result achieved, and describe the applications and significance of the results. Highlight anything novel that was discovered or improved by the thesis. Maximum length is 200 words, but try to fit into 120. Abstracts are often used for deciding if a reviewer will be suitable for the thesis; a well-written abstract thus increases the probability of getting a reviewer who will like the thesis.

Keywords: key words

Contents

Introduction	3
1 Important first chapter	5
2 Theory	7
2.1 Reinforcement learning	7
2.2 Evolutionary algorithms	7
2.3 Evolutionary strategies	8
2.3.1 CMA-ES	9
2.4 Evolutionary strategies as replacement for reinforcement learning	9
2.4.1 OpenAI Evolutionary Strategy	10
3 Results and discussion	13
Conclusion	15
Bibliography	17
A Using CoolThesisSoftware	19

Introduction

Chapter 1

Important first chapter

Chapter 2

Theory

2.1 Reinforcement learning

- problem description (state space, action space, reward...)
- various methods
 - value function
 - criterion of optimality
 - direct policy search (and various methods)

2.2 Evolutionary algorithms

Evolutionary algorithms (EA) are a type of optimisation metaheuristics inspired by the process of biological evolution. At first a number of possible solutions to the problem at hand is generated (*population*) and each solution (*individual*) is encoded (via a domain-specific encoding) and evaluated giving us the value of its *fitness*. Fitness is a function describing how good that particular individual is and it is everything that is needed for creation of Then a new population is created using a *crossover* (combination) of 1 or more individuals which are selected using the operator of *parental selection*. Each of the newly created individuals has a chance to be mutated via the *mutation* operator. Finally a new population is selected from *offsprings* and possibly the parents based of fitness and enters the next iteration of the EA and the following generation is chosen using *environmental selection* operator. The algorithm repeats until the stop condition is met, usually a set number of iterations or small improvement of fitness between 2 generations.

There are many variants of EAs such as genetic algorithms (most common), genetic programming, evolutionary programming, neuroevolution or evolutionary strategies that are further described in following chapter. [1] [2]

Algorithm 1 Evolutionary algorithm

```

1: initialize population  $P^0$  with  $n$  individuals
2: set  $t = 0$ 
3: repeat
4:    $Q^t = \{\}$ 
5:   for  $i \in \{1 \dots m\}$  do
6:      $p_1, \dots, p_\rho = \text{ParentalSelection}(P^t)$ 
7:      $q = \text{Crossover}(p_1, \dots, p_\rho)$ 
8:      $q = \text{Mutation}(q)$  with chance  $p$ 
9:      $Q^t = q \cup Q^t$ 
10:  end for
11:   $P^{t+1} = \text{EnvironmentalSelection}(Q^t \cup P^t)$ 
12:  increment  $t$ 
13: until stop criterion fulfilled

```

2.3 Evolutionary strategies

Evolutionary strategies (ES) are a type of optimisation metaheuristic which further specialises EA and restricts their level of freedom. The selection for crossover is unbiased, mutation is parametrised and thus controllable, individuals which should be put to next generation are chosen ordinaly based on fitness and individuals contain not only the problem solution but also control parameters.

More formally ES $(\mu/\rho, \kappa, \lambda)$ has μ individuals in each generation, which produces λ offsprings, each created by crossover of ρ individuals and each individual is able to survive for up to κ generations as described in algorithm 2. This notation further generalizes the old (μ, λ) and $(\mu + \lambda)$ notations, where the " " notation means $\kappa = 1$ and "+" notation $\kappa = \infty$.

To design an ES one must first select an appropriate representation for an individual and the most natural one is preferred in most cases, if all parameters are of one type (e.g. a real number) a simple vector will suffice, if the types are mixed, a tuple of vectors is required. This however causes an increased complexity of the variation operator.

As for design of the variation operator there are some guidelines that should be followed when designing it.

Algorithm 2 $(\mu/\rho, \kappa, \lambda)$ -ES

```
1: initialize population  $P^0$  with  $\mu$  individuals
2: set age for each  $p \in P^0$  to 1
3: set  $t = 0$ 
4: repeat
5:    $Q^t = \{\}$ 
6:   for  $i \in \{1 \dots \lambda\}$  do
7:     select  $\rho$  parents  $p_1, \dots, p_\rho \in P^t$  uniformly at random
8:      $q = \text{variation}(p_1, \dots, p_\rho)$  with age 0
9:      $Q^t = q \cup Q^t$ 
10:  end for
11:   $P^{t+1} = \text{select } \mu \text{ best (wrt. fitness) individuals from } Q^t \cup \{p \in P^t : \text{age}(p) < \kappa\}$ 
12:  increment age by 1 for each  $p \in P^{t+1}$ 
13:  increment  $t$ 
14: until stop criterion fulfilled
```

Reachability every solution should be reachable from any other solution in a finite number of applications of the variation operator with probability $p > 0$

Unbiasedness the operator should not favour any particular subset of solution unless provided with information about problem at hand

Control the operator should be parametrised in such way that the size of the distribution can be controlled (practice had shown that decreasing it as the optimal solution is being approached is necessary)

[3] [1]

kovariance

2.3.1 CMA-ES

TODO [4]

2.4 Evolutionary strategies as replacement for reinforcement learning

Black-box optimisation is an alternative approach to solving RL tasks also known as Direct policy search or neurevolution when applied to neural networks. It has

several attractive properties such as indifference to distribution of rewards, no need for backpropagation and tolerance of arbitrarily long episodes.

2.4.1 OpenAI Evolutionary Strategy

Compared to reinforcement learning using evolutionary strategies have the advantage of not needing a gradient of the policy performance. Also as the state transition function is not known the gradient can't be computed using backpropagation-like algorithm. Thus some noise needs to be added to make the problem smooth and the gradient to be estimable. Here is where reinforcement learning and evolutionary strategies differ, reinforcement learning adds noise in the action space (actions are chosen from a distribution) while evolutionary strategies add noise in the parameter space (parameters perturbed while actions are deterministic).

Not requiring backpropagation has several advantages over other RL methods. First the amount of computation necessary for one episode of ES is much lower (about one third, potentially even less for memory usage). Not calculating gradient using analytic methods also protects these methods from suffering from *exploding gradient* which is a common issue with recurrent neural networks. And last, the network can contain elements that are not differentiable such as hard attention.

As the ES could be seen as method for computing a derivative estimate using finite differences in randomly chosen direction it would suggest that it would scale poorly with dimensions of parameters θ same as the finite differences method. In theory the number of necessary optimisation steps should scale linearly with the dimension. That however doesn't mean that larger networks optimised using ES will perform worse than smaller ones, that depends on the difficulty (intrinsic dimension) of the problem. The network will perform the same however it will take more optimisation steps to do so.

In practice ES performs slightly better on larger networks and it is hypothesised that it is for the same reason as why it is easier to optimise large networks using standard gradient based methods: larger networks have fewer local minima.

Due to perturbing the parameters and not the actions ES are invariant to the frequency at which the agent acts in the environment. Traditional MDP-based reinforcement learning methods rely on *frameskip* as one their parameters that is crucial to get right for the optimization to be successful. While this is solvable for problems that do not require long term planning and actions, long term strategic behaviour poses a challenge and reinforcement learning needs hierarchy to be successful unlike evolutionary strategy.

- Evolution Strategies as a Scalable Alternative to Reinforcement Learning [5]
 - algorithm description
 - comparison with RL
 - parallelization
- Improving Exploration in Evolution Strategies for Deep Reinforcement Learning via a Population of Novelty-Seeking Agents [6]
 - directed exploration
 - novelty search
 - algorithms description
 - ratio of fitness and novelty and its effects

Chapter 3

Results and discussion

Conclusion

Bibliography

- [1] Günter Rudolph. “Evolutionary Strategies”. In: *Handbook of Natural Computing*. Ed. by Grzegorz Rozenberg, Thomas Bäck, and Joost N. Kok. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 673–698. ISBN: 978-3-540-92910-9. doi: 10.1007/978-3-540-92910-9_22. URL: https://doi.org/10.1007/978-3-540-92910-9_22.
- [2] P. A. Vikhar. “Evolutionary algorithms: A critical review and its future prospects”. In: *2016 International Conference on Global Trends in Signal Processing, Information Computing and Communication (ICGTSPICC)*. 2016, pp. 261–265. doi: 10.1109/ICGTSPICC.2016.7955308.
- [3] Hans-Paul Schwefel and Günter Rudolph. “Contemporary evolution strategies”. In: *Advances in Artificial Life*. Ed. by Federico Morán et al. Berlin, Heidelberg: Springer Berlin Heidelberg, 1995, pp. 891–907. ISBN: 978-3-540-49286-3.
- [4] Nikolaus Hansen. *The CMA Evolution Strategy: A Comparing Review*. 2006.
- [5] Tim Salimans et al. *Evolution Strategies as a Scalable Alternative to Reinforcement Learning*. 2017. arXiv: 1703.03864 [stat.ML].
- [6] Edoardo Conti et al. *Improving Exploration in Evolution Strategies for Deep Reinforcement Learning via a Population of Novelty-Seeking Agents*. 2018. arXiv: 1712.06560 [cs.AI].

Appendix A

Using CoolThesisSoftware

