

ALMA MATER STUDIORUM · UNIVERSITY OF BOLOGNA

Flatland Challenge



FLATLAND

Deep learning course final project

Leonardo Calbi (leonardo.calbi@studio.unibo.it)

Alessio Falai (alessio.falai@studio.unibo.it)

December 4, 2020

Contents

| | | |
|----------|------------------------------|-----------|
| 1 | Introduction | 4 |
| 2 | Background | 5 |
| 3 | Environment | 6 |
| 3.1 | Railway encoding | 6 |
| 3.2 | Observations | 6 |
| 3.2.1 | Tree | 6 |
| 3.2.2 | Binary tree | 6 |
| 3.3 | Predictions | 6 |
| 3.3.1 | Shortest path | 6 |
| 3.4 | Choices | 6 |
| 3.5 | Rewards shaping | 6 |
| 4 | Policy | 7 |
| 4.1 | Action masking | 7 |
| 4.2 | Action selection | 7 |
| 4.2.1 | ϵ -greedy | 7 |
| 4.2.2 | Boltzmann | 7 |
| 4.3 | Replay buffers | 7 |
| 4.3.1 | Uniform | 7 |
| 4.3.2 | Prioritized | 7 |
| 5 | DQN | 8 |
| 5.1 | Architectures | 8 |
| 5.1.1 | Vanilla | 8 |
| 5.1.2 | Double | 8 |
| 5.1.3 | Dueling | 8 |
| 5.2 | Bellman equation | 8 |
| 5.2.1 | Max | 8 |
| 5.2.2 | Softmax | 8 |
| 6 | GNN | 9 |
| 7 | Results | 10 |
| 8 | Conclusions | 11 |

List of Figures

Foreword

Introduction

Background

Environment

3.1 Railway encoding

3.2 Observations

3.2.1 Tree

3.2.2 Binary tree

3.3 Predictions

3.3.1 Shortest path

3.4 Choices

3.5 Rewards shaping

Policy

4.1 Action masking

4.2 Action selection

4.2.1 ϵ -greedy

4.2.2 Boltzmann

4.3 Replay buffers

4.3.1 Uniform

4.3.2 Prioritized

DQN

5.1 Architectures

5.1.1 Vanilla

5.1.2 Double

5.1.3 Dueling

5.2 Bellman equation

5.2.1 Max

5.2.2 Softmax

GNN

Results

Conclusions