Follow your Nose: Using General Value Functions for Directed Exploration in Reinforcement Learning

CCS Research

Durgesh Kalwar¹, Omkar Shelke¹, Somjit Nath¹, Hardik Meisheri¹, Harshad Khadilkar^{1,2}

> ¹TCS Research, Mumbai ²IIT Bombay

Contribution

• We extend upon temporally extended version of the ϵ -greedy exploration strategy by using auxiliary task learning with the help of General Value Functions (GVF) to perform directed exploration thereby further improving state space coverage during exploration.

Environments

Subgoal TwoRooms

MiniGrid DoorKey



• This is generalized formulation to include domain knowledge about the environment by providing GVF cumulant which also improves latent representation.



Figure 1: High-level architecture of the DEZ-greedy strategy.





Results & Discussion



Pseudo-code

Function $DEZGreedy(\epsilon, Z_{max})$:

Countdown timer $z \leftarrow 0$

```
Uniformly sampled random action w \leftarrow -1
```

Selected GVF index $g \leftarrow 0$

while True do

Observe State s

if Z == 0 then

```
if random() < \epsilon then
 Sample Duration: z \sim [1, Z_{max}]
                                               // Explore
 Sample GVF: g \sim [0, M]
if g == 0 then
      Sample action: w \leftarrow U(A)
      \mathfrak{a} \leftarrow \mathfrak{W}
 a \leftarrow argmax(Q_a^{GVF})
```



References

Will Dabney et al, Temporally-Extended ϵ -Greedy Exploration, ICLR 2021, |1| https://openreview.net/forum?id=ONBPHFZ7zG4

AAMAS, 29 May – 02 June 2023, London