

# Multi-Agent Area Coverage Control using Reinforcement Learning

Adekunle Adepegba, Suruz Miah, Davide Spinello  
Presented by Simon Hu

## PROBLEM STATEMENT

Consider a group of  $N$  homogeneous agents moving in a compact environment  $\Omega \subset \mathbb{R}^2$  where the dynamics of the agent are given by  $\dot{p}_i = g(p_i, u_i)$  where  $p_i = (x_i, y_i) \in \mathbb{R}^2$  is the agent position and  $u_i = (u_{x_i}, u_{y_i})$  represents the control input. The goal is to find a configuration of agent positions  $p = (p_1, p_2, \dots, p_N)$  such that the cost index

$$\mathcal{H}(p, t) = \int_{\Omega} \max_{i=1, \dots, N} f_i(\|p_i - q\|) \phi(q, t) \, dq$$

is maximized.

## VORONOI PARTITIONS

The Voronoi partition of  $\Omega$  is given by  $\mathcal{V} = \bigcup \mathcal{V}_i$  where each  $\mathcal{V}_i$  is given by

$$\left\{ q \in \Omega \mid \|q - p_i\| \leq \|q - p_j\|, \forall j \neq i \right\}.$$

The mass  $m_{\mathcal{V}_i}$  and center of mass  $c_{\mathcal{V}_i}$  of each  $V_i$  are given by

$$m_{\mathcal{V}_i} = \int_{\mathcal{V}_i} \phi \, dq, c_{\mathcal{V}_i} = \frac{1}{m_{\mathcal{V}_i}} \int_{\mathcal{V}_i} q \phi \, dq$$

Then the cost index  $\mathcal{H}$  can be rewritten as

$$\mathcal{H}(p, t) = \sum_{i=1}^n \int_{\mathcal{V}_i} f_i(\|p_i - q\|) \phi(q, t) \, dq.$$

[Cortés et al., 2004] showed that the optimal partition of  $\Omega$  that maximizes  $\mathcal{H}$  is the centroidal voronoi partition.

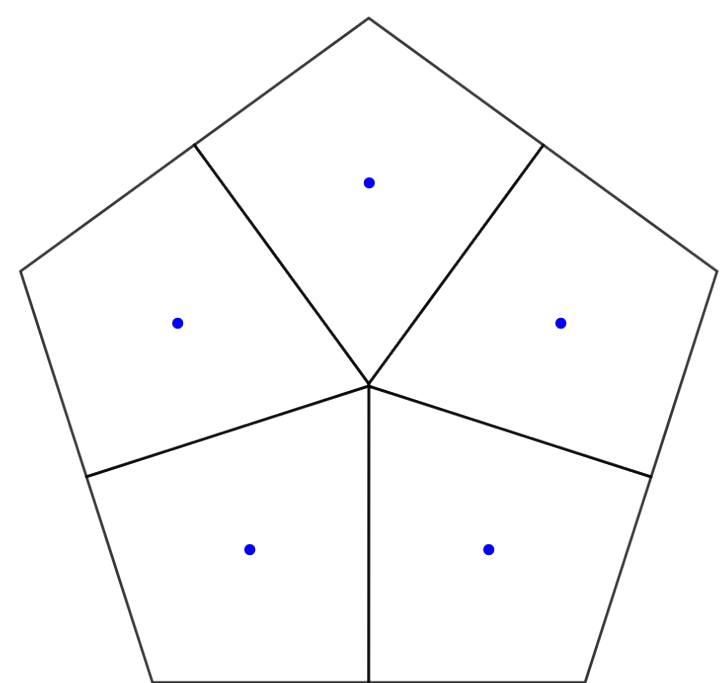


Figure 1: Example of a centroidal voronoi partition, for  $\phi(q, t) = 1$ .

## TD3 ALGORITHM

### Algorithm 1 Twin-Delayed Actor-Critic DDPG

- 1: Initialize critic networks  $Q_{\theta_1}, Q_{\theta_2}$  and actor network  $\pi_{\phi}$  with random parameters  $\theta_1, \theta_2, \phi$ .
- 2: Initialize the target networks  $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2, \phi' \leftarrow \phi$ .
- 3: Initialize a replay buffer  $R$ .
- 4: **for**  $t = 1$  **to**  $T$  **do**
- 5:   Select action with exploration noise  $a \sim \pi(s) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma)$  and record the reward  $r$  and new state  $s'$ .
- 6:   Store the tuple  $(s, a, r, s')$  into  $R$ .
- 7:   Sample minibatches of  $N$  transitions  $(s, a, r, s')$  from  $R$ .
- 8:   Smooth the target policy according to  $\tilde{a} \leftarrow \pi_{\phi'}(s) + \epsilon, \epsilon \sim \text{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$ .
- 9:   Perform double  $Q$ -learning and clip the results according to  $y \leftarrow r + \gamma \min_{1,2} Q_{\theta'_i}(s', \tilde{a})$ .
- 10:   Update the critics according to  $\theta_i \leftarrow \min_{\theta_i} N^{-1} \sum (y - Q_{\theta_i}(s, a))^2$ .
- 11:   **if**  $t \bmod d$  **then**
- 12:     Update  $\phi$  by the deterministic policy gradient  $\nabla_{\phi} J(\phi) = N^{-1} \sum \nabla_a Q_{\theta_1}(s, a) \big|_{a=\pi_{\phi}(s)} \nabla_{\phi} \pi_{\phi}(s)$ .
- 13:     Update the target networks according to  $\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i, \phi' \leftarrow \tau \phi + (1 - \tau) \phi'$ .
- 14:   **end if**
- 15: **end for**

The TD3 algorithm is an improvement on the SACDDPG algorithm, which is more vanilla, and prevents overestimation of the value function by decoupling the action selection and  $Q$ -value update.

TD3 reduces variance by updating the policy at a lower frequency than the  $Q$ -function updates.

A regularization strategy is introduced by adding a small amount of clipped random Gaussian noise to the selected action and then averages it over minibatches.

## EXTENSIONS

Placeholder.

## EXPERIMENTAL SETUP

Placeholder.

## RESULTS

Placeholder.

## References

- [1] C. Nowzari and J. Cortés, “Self-triggered coordination of robotic networks for optimal deployment,” *Automatica*, vol. 48, pp. 1077–1087, June 2012.
- [2] J. Cortés, S. Martínez, T. Karatas, and F. Bullo, “Coverage control for mobile sensing networks,” *IEEE Transactions on Robotics and Automation*, vol. 20, pp. 243–255, April 2004.
- [3] A. Adepegba, S. Miah, and D. Spinello, “Multi-agent area coverage control using reinforcement learning,” 2016.
- [4] “Deep deterministic policy gradient,” *Deep Deterministic Policy Gradient - Spinning Up Documentation*.