
Algorithm: VDN Training

- 1: Initialize θ , the parameters of individual Q-networks, and θ^- the parameters of target networks.
- 2: Initialize replay buffer \mathcal{D} // $(\mathbf{o}_t, \mathbf{a}_t, r_t, done, \mathbf{o}_{t+1})$
- 3: **while** $t < T$ **do**
- 4: Collect observations $\{o_1^t, \dots, o_n^t\}$
- 5: **for** each agent i **do**
- 6: With probability ϵ , select random action a_i^t
- 7: otherwise select $a_i^t = \arg \max_{a_i} Q_i(o_i^t, a_i)$
- 8: **end for**
- 9: Execute joint action $\mathbf{a}^t = (a_1^t, \dots, a_n^t)$
- 10: Collect r^t , $done^t$, and \mathbf{o}^{t+1}
- 11: Store $(\mathbf{o}^t, \mathbf{a}^t, r^t, done^t, \mathbf{o}^{t+1})$ in \mathcal{D}
- 12: **if** t is a training step **then**
- 13: Sample batch $\mathcal{B} = \{\mathbf{o}^b, \mathbf{a}^b, r^b, done^b, \mathbf{o}'^b\}$
- 14: Set the targets

$$y^b = r^b + \gamma(1 - done^b) \sum_i \max_{a'_i} Q_i(o_i^b, a'_i; \theta^-)$$

- 15: Perform a gradient descent using:

$$\mathcal{L}(\theta) = \frac{1}{|\mathcal{B}|} \sum_b \left(y^b - \sum_i Q_i(o_i^b, a_i^b; \theta) \right)^2$$

- 16: Every C steps, update $\theta^- \leftarrow \theta$
 - 17: **end if**
 - 18: **end while**
-