

1 Training process

Using neural networks for training of our bomberman agent, we decided to implement Rainbow (without distributional). For this technique, a number of improvements in deep learning is combined to form an integrated agent. First of all we started with Q-learning. The Q_π value is computed by choosing an action a for a given state s and calculating all of the future rewards R_n under a given policy π . Hereby future rewards are scaled by a discount factor γ to increase the weight of the immediate rewards.

$$Q_\pi(s, a) = E \left[\sum_{n=1}^{\infty} \gamma^{n-1} R_n \right] \Big|_{a, s, \pi} \quad (1)$$

The optimal Q-values are then determined by taking the maximum of all Q_π -values. To estimate these, the neural network learns a Q-function $Q(a, s, \theta_t)$ where θ_t are the weights calculated by the network. These are updated towards the target

$$Y_t^Q \equiv R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \theta_t) \quad (2)$$

according to

$$\theta_{t+1} = \theta_t + \alpha \left(Y_t^Q - Q(S_t, A_t; \theta_t) \right) \nabla_{\theta_t} Q(S_t, A_t; \theta_t) \quad (3)$$