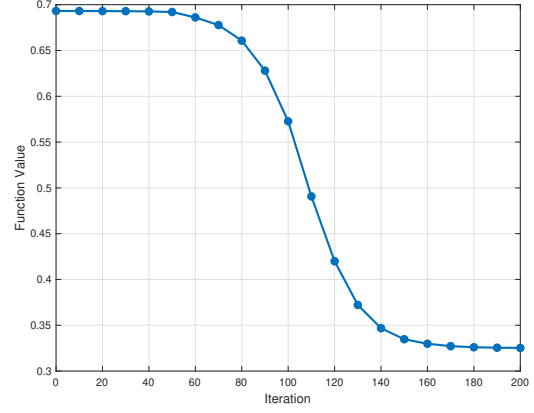


(a)



(b)

**Fig. 3.** (a) Trajectories of all five agents when initialized on the saddle point (0,0). Note that all trajectories overlap with each other, implying perfect consensus among the agents. (b) Global Function Value Dynamics. In our experiment with explicit saddle points, we consider a simple  $\{0, 1\}$ -classification task using a neural network with a single linear hidden layer and a logistic activation function. We use the cross-entropy loss function to train the network. We denote the feature vector as  $\mathbf{h} \in \mathbb{R}^M$  and the binary class label as  $y \in \{-1, 1\}$ . For the fully connected hidden layer, we represent the weights as  $\mathbf{W}_2 \in \mathbb{R}^{L \times M}$  and  $\mathbf{W}_1 \in \mathbb{R}^L$ . The output is of the form:  $\hat{y} = \frac{1}{1 + e^{-\langle \mathbf{h}, \mathbf{W}_2^\top \mathbf{W}_1 \rangle}}$ . Under the commonly used cross-entropy loss

function, the objective function is of the following form:  $L(\mathbf{W}_1, \mathbf{W}_2) = \log(1 + e^{-\langle \mathbf{h}, \mathbf{W}_2^\top \mathbf{W}_1 \rangle})$ . To visualize the evolution of optimization variables under our algorithm, we considered the scalar case with  $L = M = 1$  and purposely initialized all the agents from the strict saddle point (0, 0). Synthetic data samples ( $\mathbf{h}$  and  $y$ ) are randomly generated under the constraint  $\mathbf{h}y = 1$  and the regularization parameter is set to  $\rho = 0.1$ . For our algorithm, the stepsize, noise amplitude, noise-injection interval and clipping threshold were set to  $\alpha = 0.1$ ,  $\theta = 0.1$ ,  $\mathcal{K}_0 = 50$  and  $c_0 = 1.0$ , respectively. It can be seen that due to the noise and gradient clipping effect, all five agents collectively move along the descending direction, clearly indicating that our algorithm can ensure efficient escape from saddle points.