Figure 4: Results of different methods when using different step size schemes. Column (a): constant $\epsilon_t = \epsilon_0$; Column (b): decaying step size $\epsilon_t = \epsilon_0(b + t)^{-\gamma}$; Column (c): Adagrad (whose master step size we denote by $\epsilon_0$). We search the best $\epsilon_0$ in the grid $[1e-3, 1e-4, 1e-5, 1e-6]$ that achieves the lowest constraint loss at the end of the training. For the other parameters, we use fixed $\gamma = 0.55$ and $b = 1$ for decaying step size, and the default parameters of Adagrad in PyTorch except the master step size $\epsilon_0$.