

V této práci zkoumáme možnost tvorby nepřátelských vzorků pomocí generativních modelů. Použijeme generativní modely k vytvoření nepřátelských vzorků pomocí perturbace latentních vektorů a využití některých vlastností klasifikátoru. Součástí práce je návrh dvou algoritmů. První algoritmus využívá lineární interpolace v kombinaci s bisekcí k získání vzorků z rozhodovací hranice klasifikátoru. Druhý využívá gradient k vytvoření potenciálně nepřátelských vzorků, podobně jako algoritmus FGSM. Tyto algoritmy použijeme na datasety MNIST, SVHT a CIFAR a vytvoříme sadu potenciálně nepřátelských vzorků a ukážeme, že v nich existují vzorky, které člověk klasifikuje jinak než stroj. Dále ukážeme, že náš typ útoku dokáže obejít běžně používané obrany.