

Hluboké neuronové sítě v poslední době dosahují vysoké úspěšnosti na mnoha úlohách, zejména klasifikaci obrázků. Tyto modely jsou ovšem snadno ovlivnitelné lehce pozměněnými vstupy zvanými matoucí vzory. Matoucí vzory mohou značně snižovat úspěšnost a tak ohrozit systémy, které modely strojového učení využívají. V této práci přinášíme rešerši literatury o matoucích vzorech. Dále navrhuje nové obrany proti matoucím vzorům: síť kombinující RBF jednotky s konvolucí, kterou testujeme na datové sadě MNIST a která má lepší úspěšnost než CNN trénovaná pomocí matoucích vzorů, a diskretizaci vstupního prostoru, kterou testujeme na datových sadách MNIST a ImageNet a dosahujeme slibných výsledků. Na závěr zkoumáme možnost generování matoucích vzorů bez přístupu ke vstupu, který má být pozměněn.