



ČESKÉ VYSOKÉ UČENÍ TECHNICKÉ V PRAZE
Fakulta jaderná a fyzikálně inženýrská



Robustní strojové učení a adversariální vzorky

Robust machine learning and adversarial examples

Bakalářská práce

Autor: **Pavel Jakš**
Vedoucí práce: **Mgr. Lukáš Adam, Ph.D.**
Akademický rok: 2021/2022

- Zadání práce -

- Zadání práce (zadní strana) -

Poděkování:

Chtěl bych zde poděkovat především svému školiteli - panu doktoru Adamovi - za pečlivost, ochotu, vstřícnost a odborné i lidské zázemí při vedení mé bakalářské práce.

Čestné prohlášení:

Prohlašuji, že jsem tuto práci vypracoval samostatně a uvedl jsem všechnu použitou literaturu.

V Praze dne 7. července 2022

Pavel Jakš

Robustní strojové učení a adversariální vzorky

Obor: Matematická informatika

Vedoucí práce: Mgr. Lukáš Adam, Ph.D., Katedra počítačů, Fakulta elektrotechnická, České vysoké učení technické v Praze, Karlovo náměstí 13, 121 35, Praha 2

[illegible]

Klíčová slova: klíčová slova (nebo výrazy) seřazená podle abecedy a oddělená čárkou

Robust machine learning and adversarial examples

[illegible]

Key words: keywords in alphabetical order separated by commas

Obsah

Úvod	11
1 Neuronové sítě	13
1.1 Hluboká dopředná neuronová síť	13
1.2 Konvoluční síť	14
2 Učení neuronové sítě	17
2.1 Účelové funkce	17
2.2 Algoritmus zpětného šíření chyby	18
2.3 Algoritmy učení	18
2.3.1 Gradientní sestup	18
2.3.2 Stochastický gradientní sestup	18
2.3.3 Adam	18
3 Adversariální vzorky	19
3.1 Metody generování adversariálních vzorků	19
3.1.1 FGSM	19
3.1.2 Iterativní FGSM	19
4 Robustní učení neuronové sítě	21
Závěr	23

Úvod

Pojem neuronové sítě představuje výpočetní jednotku, která svou univerzálností nachází uplatnění v mnoha disciplínách.

Kapitola 1

Neuronové sítě

Princip fungování neuronové sítě spočívá v poskládání celku z dílčích výpočetních jednotek - umělých neuronů. Takovýto neuron je standardně funkcí více proměnných, jehož výstup je proměnná jediná. Typickým modelem umělého neuronu je funkce $f : \mathbb{R}^n \rightarrow \mathbb{R}$ definovaná předpisem

$$f(x_1, \dots, x_n) = \sigma\left(\sum_{i=1}^n w_i x_i + b\right), \quad (1.1)$$

kde n je počet vstupujících proměnných, w_i jsou tzv. váhy (w z anglického slova weight), b je práh (b z anglického slova bias), σ označuje tzv. aktivační funkci.

Roli vstupujících proměnných mohou hrát např. hodnoty RGB pixelů barevných obrázků, je-li aplikační klasifikace obrázků, nebo výstupy jiných neuronů. Pod pojmem váha se skrývá míra ovlivnění výstupu neuronu daným vstupem. Je-li váha u nějakého vstupu vysoká, pak je výstup citlivější na daný vstup. Prah pro změnu určuje posunutí citlivosti neuronu na všechny vstupy jako celku.

Poslední, avšak velmi důležitou charakteristikou tohoto modelu neuronu je aktivační funkce. Za aktivační funkci lze vzít libovolnou funkci $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, existuje však základní sada:

- Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$
- ReLU: $\sigma(x) = \max(0, x)$
- LeakyReLU: $\sigma(x) = \max(0, x) + \alpha * \min(x, 0)$, kde $\alpha \in \mathbb{R}^+$
- Tanh: $\sigma(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$

Tyto funkce lze doplnit o jejich mírné modifikace. Moderní doporučenou praxí je užívat ReLU jako aktivační funkci a (1.1) jako model neuronu dle [1].

1.1 Hluboká dopředná neuronová síť

Je-li pojem umělého neuronu objasněn, lze se přesunout k jeho užití v neuronových sítích. Základní myšlenkou těchto sítí je vhodné poskládání umělých neuronů do vrstev, které dohromady tvoří síť neuronů. Taková vrstva je potom trojího druhu - vstupní, výstupní a skrytá. *Vstupní vrstva* je množina umělých neuronů, které mají za vstup výstupy problému, jehož je neuronová síť řešením. Za vstup si lze představit matici černobílých pixelů, které představují obrázek číslice, kterou je cíl klasifikovat. *Výstupní vrstva* sestává z neuronů, které mají za vstup výstupy neuronů předchozí vrstvy. Výstupem této vrstvy pak bude řešení daného problému - například klasifikace číslice. Posledním druhem vrstvy je *vrstva*

skrytá. Takováto vrstva má za vstupy výstupy vrstvy předcházející a její výstupy slouží jako vstupy pro vrstvu nadcházející. Má-li neuronová síť tuto architekturu, hovoří se o *dopředné neuronové síti*. Má-li navíc alespoň jednu skrytou vrstvu, lze mluvit o *hluboké dopředné neuronové síti*.

Se znalostí pojmu vrstvy neuronů lze přistoupit k poznámce o tzv. *softmax funkci*. Jedná se o vektorovou funkci $s : \mathbb{R}^m \rightarrow \mathbb{R}^m$, kde

$$s(x_1, \dots, x_m)_i = \frac{e^{x_i}}{\sum_{j=1}^m e^{x_j}},$$

kde $i \in \hat{m}$. Její užití je nasnadě: Výstup této funkce lze totiž interpretovat jako diskretní pravděpodobnostní distribuci, a proto ji lze užít jako aktivační funkci výstupní vrstvy, je-li cílem dané neuronové sítě klasifikace vstupu do kategorií.

Další poznámka se bude věnovat zjednodušení zápisu akce vrstvy na vstup. Podle modelu neuronu v (1.1) se akce jednoduho neuronu na vstup sestává z násobení, následného sčítání, přičtení prahu a aplikací aktivační funkce. Tato procedura nastává pro každý neuron ve vrstvě. Tak lze sestavit z jednotlivých vah $w_i^{(j)}$ (i -tá váha j -tého neuronu ve vrstvě) matici \mathbb{A} , jejímiž prvky jsou právě ony váhy $(\mathbb{A})_{j,i} = w_i^{(j)}$, z prahů pak vektor b , jehož j -tá složka je rovna prahu j -tého neuronu. Dále zaved' me vektorovou funkci $s : \mathbb{R}^m \rightarrow \mathbb{R}^m$ - at' už jako výše zmíněnou softmax funkci, nebo jako po složkách aplikovanou libovolnou aktivační funkci σ ve smyslu $s(x_1, \dots, x_m)_i = \sigma(x_i)$ pro $i \in \hat{m}$. Pak lze psát, že aplikace vrstvy neuronů je zobrazení $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ působící na vektor x následovně:

$$\phi(x) = s(\mathbb{A}x + b) \quad (1.2)$$

Tedy stěžejní operací se stává maticové násobení, respektive násobení vektoru maticí zprava.

Při tomto si lze povšimnout, že takováto neuronová síť má řadu parametrů, o kterých není jasné jak je správně nastavit. Některé parametry (například váhy a prahy) se nastavují během učení neuronové sítě, čemuž je věnována samostatná kapitola. Potom tu jsou parametry, jejichž charakter je poněkud odlišný. Jedná se o ty parametry, které zůstávají během života neuronové sítě netknuté. Jako příklad lze uvést počet neuronů ve skryté vrstvě, který se promítne v rozměrech matice vah či dimenzionalitě výstupu vrstvy. Takovýmto parametrům je přisuzován název hyper-parametry.

1.2 Konvoluční síť

Konvoluční síť nebo též *konvoluční neuronové síť* přinášejí svou architekturou nové možnosti zpracování dat se specifickou strukturou, do které patří například časové řady, obrázky nebo videa. Středobodem konvolučních sítí je, jak již název napovídá, operace *konvoluce*. Ta nahrazuje maticové násobení, kterým lze reprezentovat operace ve výše popsaném modelu hluboké dopředné sítě.

Operace *konvoluce* je ve vší obecnosti operace mezi dvěma číselnými funkcemi g a h se stejným definičním oborem, jejíž výstupem je nová číselná funkce standardně označovaná jako $g * h$. Uved' me zde definici konvoluce pro reálné funkce definované na \mathbb{R}^d , tedy $g, h : \mathbb{R}^d \rightarrow \mathbb{R}$:

$$(g * h)(t) = \int_{\mathbb{R}^d} g(x)h(t - x)dx$$

Důležitým předpokladem pro možnost konvoluce je samozřejmě konvergence integrálu na pravé straně.

Ačkoliv je konvoluce komutativní operací, v kontextu strojového učení se mezi oběma funkcemi vstupujícími do konvoluce rozlišuje. Funkce vstupující jako první se nazývá vstup a druhá funkce se nazývá jádrem. Dále se v kontextu konvolučních sítí standardně objevují diskretní funkce, které nabývají nenulových hodnot pouze v konečně mnoha bodech. Potom integrál přes \mathbb{R}^d přechází v konečnou sumu:

$$(g * h)(i_1, \dots, i_d) = \sum_{j_1} \dots \sum_{j_d} g(j_1, \dots, j_d)h(i_1 - j_1, \dots, i_d - j_d) \quad (1.3)$$

Díky komutativitě konvoluce lze též psát:

$$(g * h)(i_1, \dots, i_d) = \sum_{j_1} \dots \sum_{j_d} g(i_1 - j_1, \dots, i_d - j_d) h(j_1, \dots, j_d) \quad (1.4)$$

Při aplikaci komutativity došlo k tzv. *překlopení jádra* (termín pochází z anglického kernel flipping). Za vynechání překlopení jádra lze dojít ke *křížové korelaci*:

$$(g * h)(i_1, \dots, i_d) = \sum_{j_1} \dots \sum_{j_d} g(i_1 + j_1, \dots, i_d + j_d) h(j_1, \dots, j_d) \quad (1.5)$$

Mnoho knihoven zabývajících se neuronovými sítěmi dle [1] implementují křížovou korelaci namísto konvoluce, ačkoliv tuto svou implementaci nazývají konvolucí.

Další nedílnou součástí konvolučních sítí je tzv. *pooling*. Spolu s konvolucí tvoří mocný nástroj, který ve formě konvolučních a pooling vrstev hlubokých neuronových sítí přináší například invarianci sítě vůči malému posunutí vstupu (dle [1]).

Pooling je funkce, která nahrazuje hodnoty v bodech nějakou souhrnou statistikou určitého okolí daného bodu. Např. *max pooling* aplikovaný na matici se podívá na obdélníkové okolí předem definovaných rozměrů daného bodu a jako svůj výstup vybere maximální hodnotu nalezenou v onom okolí. Jiné oblíbené pooling funkce zahrnují funkce reportující průměr či L^2 normu daného obdelníkového okolí.

Standardní konvoluční vrstva neuronové sítě pak sestává ze tří fází. První fáze provádí paralelně několik konvolucí, které produkují sadu aktivací. Druhá fáze, někdy označovaná jako *detekční fáze*, aplikuje na výstupy první fáze aktivační funkci. Třetí fáze potom provádí *pooling*.

Kapitola 2

Učení neuronové sítě

Předchozí kapitola představuje neuronové sítě jakožto složené zobrazení s mnoha parametry. Aby takováto neuronová síť byla k něčemu užitečná, např. ke klasifikaci obrázků, musí dané zkonstruované zobrazení vracet smysluplné výsledky k daným vstupům. Toho se v praxi docílí vhodným nastavením hyper-parametrů sítě a následným nalezením hodnot parametrů daného složeného zobrazení, které odpovídají funkční neuronové síti. Toto hledání parametrů se též nazývá jako *učení neuronové sítě* a provádí se metodami numerické optimalizace jistého vhodně zvoleného kritéria, které se označuje jako *účelová* či *ztrátová funkce*.

2.1 Účelové funkce

Nutným předpokladem ke konstrukci vhodné účelové funkce je tzv. *trénovací sada* vzorků a k nim příslušné *značky*. Jedná se o množinu možných vstupů, které jsou vybaveny správným výstupem. Je-li dána trénovací sada a značky, lze definovat účelovou funkci jako jakési měřidlo ukazující, jak moc se trénovaná neuronová síť mýlí, je-li vpuštěna na vzorky trénovací sady. Tomuto přístupu k učení neuronové sítě se také říká učení s učitelem.

Jedna z klasických účelových funkcí je funkce střední kvadratické chyby. Je dána přepisem:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m (F_{\theta}(x^{(i)})_j - y_j^{(i)})^2, \quad (2.1)$$

nebo

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N \|F_{\theta}(x^{(i)}) - y^{(i)}\|_2^2 \quad (2.2)$$

kde $x^{(i)}$ je i -tý vektor trénovací sady, $y^{(i)}$ je i -tý vektor trénovacích značek, F_{θ} neuronová síť jakožto funkce $F_{\theta} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ parametrizovaná parametry θ .

Další účelová funkce, která nachází uplatnění v klasifikačních problémech, je křížová entropie. Lze ji vypočítat následovně:

$$L(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m y_j^{(i)} * \ln(F_{\theta}(x^{(i)})_j) \quad (2.3)$$

nebo

$$L(\theta) = -\frac{1}{N} \sum_{i=1}^N H(y^{(i)}, F_{\theta}(x^{(i)})), \quad (2.4)$$

kde H označuje právě onu křížovou entropii mezi pravděpodobnostními distribucemi. Připomeňme, že klasifikační neuronová síť produkuje diskrétní pravděpodobnostní distribuce, a proto lze na výstup takovéto neuronové sítě a její značky (také pravděpodobnostní distribuce) aplikovat křížovou entropii.

2.2 Algoritmus zpětného šíření chyby

Nejčastější metody učení neuronové sítě ve svém chodu pracují s gradientem účelové funkce. K tomu je tedy nutné umět onen gradient spočítat. A k tomu právě slouží *algoritmus zpětného šíření chyby* (angl. *backpropagation*).

2.3 Algoritmy učení

2.3.1 Gradientní sestup

2.3.2 Stochastický gradientní sestup

2.3.3 Adam

Kapitola 3

Adversariální vzorky

3.1 Metody generování adversariálních vzorků

3.1.1 FGSM

3.1.2 Iterativní FGSM

Kapitola 4

Robustní učení neuronové sítě

Závěr

Text závěru....

Literatura

- [1] I. Goodfellow, Y. Bengio, A. Courville, *Deep Learning*. MIT Press, 2016.
- [2] I. Goodfellow, J. Shlens, C. Szegedy, *Explaining and Harnessing Adversarial Examples*. In 'International Conference on Learning Representations', ICLR 2015.
- [3] J. Nocedal, S. Wright, *Numerical optimization*. Springer Science & Business Media, 2006.
- [4] M. A. Nielsen, *Neural Networks and Deep Learning*. Determination Press, 2018.