

# Hipoteza "Lottery winning ticket"

1<sup>st</sup> Michał Sumiński

Politechnika Łódzka

PŁ

Łódź, Polska

249999

## I. STRESZCZENIE

Jednym z obecnych wyzwań trenowania sieci neuronowych jest tzw. overparametrization (pol. przeparametryzowanie), czyli problem z ogromną liczbą parametrów modelu. Technika 'pruningu' sieci neuronowych może znacząco zredukować liczbę tych parametrów (nawet powyżej 90%) oraz potrzebna pamięć, zwiększając przy tym wydajność oraz dokładność modelu. W niniejszej pracy zdecydowano się przetestować 3 różne warianty wykorzystania mechanizmu pruningu: pruning iteracyjny, pruning iteracyjny z przywracaniem początkowych wag (Lottery winning ticket), pruning iteracyjny z losową reinicjalizacją wag na sieciach CNN oraz GNN. Wyniki zaprezentowano w tabelach 1-5. Eksperymenty pokazują, że każdy z wariantów daje dokładność bardzo bliską oryginalnej, przy wykorzystaniu  $< 5\%$ .

## II. WPROWADZENIE

Artykuł ten poświęcony jest przebadaniu na własną rękę hipotezy Lotery winning ticket" zaproponowanej przez Jonathana Frankle'a oraz Michaela Carbin'a [1] rozszerzająca podstawowe pojęcie pruningu o dodatkowy mechanizm przywracania wag z początkowej inicjalizacji. Hipoteza ta brzmi następująco: *Losowo zainicjowana, sieć neuronowa zawiera podsieć, która jest zainicjowana w taki sposób, że podczas uczenia w izolacji może odpowiadać testowej dokładności oryginalnej sieci, kiedy jest wytrenowana przez co najwyżej taka sama liczba iteracji.* Możemy rozpatrzeć następujący przykład:

- Tworzymy sieć  $X$ , która składa się ze 100 wag
- Po wytrenowaniu tej sieci przez  $e$  epok osiągamy wartość  $accuracy(X) = a$ .
- Obcinamy wagi o połowę, zostaje ich 50, przywracamy wagi z początkowej inicjalizacji, otrzymujemy sieć  $X'$ .
- Trenujemy sieć  $X'$  przez  $e'$  epok, gdzie  $e' < e$ .
- Otrzymujemy miarę  $accuracy(X') = a'$ , gdzie zgodnie z hipotezą  $a' < a$ .

## III. POWIĄZANE PRACE

Problem przeparametryzowania jest rozwijany na wiele sposobów, Jonathan Frankle oraz Michael Carbin [1] w swojej pracy przytaczają między innymi takie prace jak:

- Restrykcja optymalizacji do małych, losowych zbiorów parametrów z całej przestrzeni parametrów. [6]
- Modele SqueezeNet, MobileNet, które są o rząd wielkości niższe od standardowych architektur. [7]

- Destylacja (ang. Distillation) – przekazanie wiedzy z większego modelu do mniejszego. [8]

## IV. METODA

### A. Konwolucyjne sieci neuronowe (CNN)

Do eksperymentów zdecydowano się użyć architektury sieci CNN z zadania nr.2, tzn. sieć CNN z dwoma warstwami konwolucyjnymi typu Conv2d oraz jedną warstwą w pełni połączoną, użyta funkcja aktywacji w warstwach konwolucyjnych to leaky\_relu, a w warstwie w pełni połączonej to sigmoid, użyto optymalizatora Adagrad oraz krzyżowej entropii jako funkcji celu. Całość została zaimplementowana przy użyciu biblioteki PyTorch[2] (a ściślej jej podzbiorem torchvision[4]). Użyty został zbiór danych MNIST.

Zdecydowano się przeprowadzić 3 różne eksperymenty, tj. pierwszy - pruning iteracyjny, tzn. uczenie modelu przez 'x' iteracji, następnie wykonanie operacji pruningu, i tak kilka razy, drugi - również pruning iteracyjny natomiast po każdej operacji pruningu pozostałe wagi zostawały losowo reinicjalizowane, trzeci - również pruning iteracyjny, a po każdej operacji pruningu zostały przywracane wagi z pierwszej inicjalizacji - jest to implementacja właśnie tytułowego "Winning ticketa". Warto przypomnieć, że oryginalnie model po 10 epokach osiągnął wartość miary  $accuracy$  równą **(0.988)**.

Do pruningu zdecydowano się zastosować również bibliotekę PyTorch. Pruning za każdym razem wykonywany był globalnie (global pruning) na 2 warstwach konwolucyjnych, tzn. że gdy podawana była wartość parametru pruningu, np. 0.2 mówiła o tym jaki procent ( $0.2 = 20\%$ ) wag ma zostać usunięty, to nie dotyczyła ona każdej warstwy pojedynczo, tzn. z każdej warstwy nie zostało usuwane 20% wag tylko. 20% wag zostało usuwanych z obu warstw, np. z pierwszej warstwy zostało usunięte 5% wag, a z drugiej 15%.

### B. Grafowe sieci neuronowe (GNN)

Ten sam mechanizm co w powyższym rozdziale o sieciach konwolucyjnych zastosowano do sieci grafowej. Eksperymenty zostały przeprowadzone na architekturze sieci z zadania nr.2, tzn. sieć GNN z dwoma warstwami konwolucyjno-grafowymi typu GCNConv oraz jedną warstwą w pełni połączoną, użyta funkcja aktywacji w warstwach konwolucyjnych to leaky\_relu, a w warstwie w pełni połączonej to sigmoid, użyto optymalizatora Adam oraz krzyżowej entropii jako funkcji celu. Całość została zaimplementowana przy

użyciu biblioteki PyTorch[2] (a ściślej jej podzbiorowi pytorch\_geometric[3]). Użyty został zbiór danych CORA. Warto przypomnieć, że oryginalnie model po 10 epokach osiągnął wartość miary accuracy równa **(0.792)**.

### C. Trening

Treningi modeli wyglądały następująco:

- 1) zainicjalizowanie modelu (w przypadku "Winning ticket'ów" po tym kroku należało skopiować macierz wag modelu w celu późniejszego ich przywrócenia)
- 2) pierwszy trening modelu (dla sieci CNN 10 iteracji, dla sieci GNN 1000)
- 3) wykonanie operacji pruningu 20%
- 4) w przypadku eksperymentów z losową reinicjalizacją bądź "Winning ticket'ami" w tym kroku zostały przywracane bądź reinicjalizowane wagi
- 5) ponowny trening modelu
- 6) wykonanie operacji pruningu 52%
- 7) Powtórzenie kroków 4,5
- 8) wykonanie operacji pruningu 80%
- 9) Powtórzenie kroków 4,5
- 10) wykonanie operacji pruningu 96%
- 11) Powtórzenie kroków 4,5
- 12) wykonanie operacji pruningu 99.6%(GNN) bądź 99.96%(CNN)
- 13) Powtórzenie kroków 4,5

## V. WYNIKI

Poniżej przedstawiono tabele z wynikami eksperymentów dla sieci CNN, liczba epok w każdej kolejnej iteracji pruningu wynosiła 10:

TABLE I  
PRUNING ITERACYJNY

| %pruningu | accuracy |
|-----------|----------|
| 20        | 0.9884   |
| 52        | 0.9897   |
| 80        | 0.9904   |
| 96        | 0.9888   |
| 99.96     | 0.1135   |

TABLE II  
PRUNING ITERACYJNY Z ZASTOSOWANIEM "WINNING TICKET'A"

| %pruningu | accuracy |
|-----------|----------|
| 20        | 0.9889   |
| 52        | 0.9896   |
| 80        | 0.9905   |
| 96        | 0.9883   |
| 99.96     | 0.5114   |

TABLE III  
PRUNING ITERACYJNY Z LOSOWĄ REINICJALIZACJĄ WAG

| %pruningu | accuracy |
|-----------|----------|
| 20        | 0.9895   |
| 52        | 0.9898   |
| 80        | 0.9903   |
| 96        | 0.9888   |
| 99.96     | 0.4238   |

Podobnie jak dla sieci konwolucyjnych przeprowadzono dokładnie te same 3 eksperymenty dla sieci GNN. Poniżej przedstawiono tabele z wynikami eksperymentów, liczba epok w ka zdej kolejnej iteracji pruningu wynosiła 1000:

TABLE IV  
PRUNING ITERACYJNY

| %pruningu | accuracy |
|-----------|----------|
| 20        | 0.9884   |
| 52        | 0.9897   |
| 80        | 0.9904   |
| 96        | 0.9888   |
| 99.6      | 0.361    |

TABLE V  
PRUNING ITERACYJNY Z ZASTOSOWANIEM "WINNING TICKET'A"

| %pruningu | accuracy |
|-----------|----------|
| 20        | 0.9889   |
| 52        | 0.9896   |
| 80        | 0.9905   |
| 96        | 0.9883   |
| 99.6      | 0.684    |

TABLE VI  
PRUNING ITERACYJNY Z LOSOWĄ REINICJALIZACJĄ WAG

| %pruningu | accuracy |
|-----------|----------|
| 20        | 0.9895   |
| 52        | 0.9898   |
| 80        | 0.9903   |
| 96        | 0.9888   |
| 99.96     | 0.319    |

## VI. PODSUMOWANIE

W niniejszych eksperymentach hipoteza "Winning ticket'a" nie rysuje się tak klarownie jak to miało miejsce u J. Franke'a, M. Carbin'a [1]. Wyniki z każdego z 3 eksperymentów dają praktycznie taką samą wartość miary accuracy - różnica jest minimalna. Natomiast można spostrzec jedną zależność, zarówno w CNN jak i w GNN sieć z "Winning ticket'em" osiągnęła największą wartość miary accuracy przy największym procencie usuniętych wag, co można zaobserwować w ostatnich wierszach każdej z tabel. Widać zatem, że gdy pozostały procent parametrów wynosił < 1%, to sieć z "Winning ticket'em" radziła sobie najlepiej. Może to być pewien punkt zaczepienia. Być może przebadane sieci były zbyt małe,

posiadały zbyt mało parametrów, aby mechanizm "Winning ticket'a" zadziałał. Co jednak jest pewne to to, że w każdym z eksperymentów wystarczył  $< 5\%$  parametrów, aby osiągnąć wyniki porównywalne do pełnej, gęstej sieci.

## VII. SPIS LITERATURY

### REFERENCES

- [1] J. Frankle, M. Carbin, "The lottery ticket hypothesis: finding sparse, trainable neural networks".
- [2] <https://pytorch.org/>
- [3] <https://pytorch-geometric.readthedocs.io/en/latest/>
- [4] <https://pytorch.org/vision/stable/index.html>
- [5] [https://pytorch.org/tutorials/intermediate/pruning\\_tutorial.html](https://pytorch.org/tutorials/intermediate/pruning_tutorial.html)
- [6] Chunyuan Li, Heerad Farkhoor, Rosanne Liu, and Jason Yosinski. Measuring the intrinsic dimension of objective landscapes. Proceedings of ICLR, 2018.
- [7] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. arXiv preprint arXiv:1602.07360, 2016.
- [8] Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In Advances in neural information processing systems, pp. 2654–2662, 2014