

ПРУНІНГ ЗГОРТКОВИХ НЕЙРОННИХ МЕРЕЖ ЗА ДОПОМОГОЮ ІНТЕРПРЕТОВАНІСТІ МЕРЕЖ КОЛМОГОРОВА-АРНОЛЬДА

Єфанов І.С.¹, Шаповал Н.В.

Національний технічний університет України
«Київський політехнічний інститут імені Ігоря Сікорського», Київ, Україна

¹ efanov.illiya@lil.kpi.ua

Запропоновано новий метод структурного прунінгу згорткових нейронних мереж, заснований на аналізі важливості ознак через інтерпретовані мережі Колмогорова-Арнольда (KAN). Метод використовує KAN-шар як вузьке місце для ідентифікації надлишкових фільтрів у CNN. Розроблено критерій важливості на основі норми коефіцієнтів В-сплайнів та алгоритм ітеративного прунінгу з донавчанням. Експериментальна валідація на CIFAR-10 показала перевагу над класичними методами: досягнуто 94.2% точності при стисненні моделі на 34%, що на 1.4% краще за magnitude-based pruning та на 0.7% краще за knowledge distillation.

Ключові слова: прунінг нейронних мереж, мережі Колмогорова-Арнольда, стиснення моделей, інтерпретованість, глибоке навчання, В-сплайни.

1. ВСТУП

Сучасні згорткові нейронні мережі досягають високої точності у задачах комп'ютерного зору, але їхнє розгортання на пристроях з обмеженими ресурсами (Edge AI) залишається складною проблемою. Моделі, такі як ResNet-50 або VGG-16, містять десятки мільйонів параметрів та вимагають гігабайти пам'яті, що робить їх непрактичними для мобільних пристроїв, IoT систем та автономного обладнання [1].

За оцінками, тренування великих мовних моделей може коштувати мільйони доларів та споживати енергію, еквівалентну викидам CO₂ від декількох автомобілів протягом їхнього життєвого циклу [2]. Для компаній, що надають сервіси на основі штучного інтелекту, навіть невелике зменшення розміру моделі може призвести до значної економії коштів на масштабі.

Методи стиснення нейронних мереж, зокрема прунінг (видалення надлишкових параметрів), дозволяють зменшити розмір та обчислювальну складність моделей [3]. Традиційні підходи базуються на евристичних критеріях важливості, найпоширенішим з яких є величина ваг (magnitude-based pruning) [4]. Однак такі методи не враховують функціональну роль компонентів мережі та їхній семантичний внесок у фінальне рішення.

Альтернативні методи включають дистилляцію знань [5], де менша модель-учень навчається імітувати велику модель-вчителя, та квантизацію [6], що зменшує точність представлення ваг. Гіпотеза лотерейного квитка [7] показала, що щільні мережі містять розріджені підмережі, які можуть досягати порівнянної точності при навчанні з правильною ініціалізацією.

Паралельно з розвитком великих моделей зростає попит на розгортання ШІ безпосередньо на кінцевих пристроях (Edge AI) – смартфонах, IoT-пристроях, автономних транспортних засобах, медичному обладнанні [1, 3]. Такі застосування мають жорсткі

обмеження: обмежену обчислювальну потужність, малий обсяг пам'яті, обмежений час автономної роботи та вимоги до реального часу. Більше того, деякі застосування, такі як медична діагностика або критична інфраструктура, не можуть покладатися на хмарні сервіси через питання приватності, безпеки та надійності з'єднання.

Обмеження Edge-пристроїв є жорсткими:

- **Пам'ять:** Типовий смартфон має 4–8 GB RAM, але більша частина зайнята операційною системою та іншими застосунками. Для моделі може залишатися лише 100–500 MB.
- **Обчислювальна потужність:** Мобільні процесори та GPU мають продуктивність на 1–2 порядки нижчу за серверні GPU (наприклад, 1–2 TFLOPS проти 100+ TFLOPS у A100).
- **Енергоспоживання:** Інференс моделі не повинен швидко розряджати батарею.
- **Латентність:** Для багатьох застосувань (розпізнавання голосу, AR/VR, автономне водіння) час відгуку повинен бути менше 50–100 мс.

Нещодавно представлені мережі Колмогорова-Арнольда (KAN) [8] відкривають нові можливості для інтерпретації нейронних мереж. На відміну від багатошарових перцептронів (MLP), де нелінійності є фіксованими функціями активації на нейронах, KAN мають навчені активаційні функції на зв'язках, параметризовані гладкими B-сплайнами [8]. Це робить KAN значно більш інтерпретованими – кожне з'єднання можна візуалізувати та аналізувати [8].

У даній роботі пропонується використати підвищену інтерпретабельність KAN для створення нового критерію важливості при структурному прунінгу згорткових шарів. Ключова ідея полягає у використанні KAN-шару як інтерпретованого вузького місця (bottleneck), що дозволяє оцінити функціональну важливість кожного згорткового фільтра через аналіз його впливу на латентний простір.

2. ТЕОРЕТИЧНІ ОСНОВИ

2.1. Теорема Колмогорова-Арнольда

Теорема представлення Колмогорова-Арнольда (1957) стверджує, що будь-яка неперервна функція змінних може бути представлена як композиція неперервних одновимірних функцій та операції додавання [9]:

$$f(x_1, \dots, x_n) = \sum_{q=0}^{2n} \Phi_q \left(\sum_{p=1}^n \varphi_{q,p}(x_p) \right), \quad (1)$$

де $\varphi_{q,p}: [0,1] \rightarrow R$ – внутрішні функції, $\Phi_q: R \rightarrow R$ – зовнішні функції. Теорема показує, що багатовимірні залежності можуть бути розкладені на комбінацію одновимірних перетворень [9, 10]. Це частково вирішило тринадцяту проблему Гільберта про неможливість представлення функцій багатьох змінних через суперпозиції функцій меншої кількості змінних [9, 10]. Однак оригінальна теорема має обмеження для практичного застосування: внутрішні функції $\varphi_{q,p}$ є недиференційовними майже скрізь та мають фрактальну природу, що ускладнює їх апроксимацію та оптимізацію градієнтними методами [8].

2.2. B-сплайни та їх властивості

B-сплайни (базисні сплайни) є кусково-поліноміальними функціями, що забезпечують ефективний спосіб апроксимації гладких кривих [11]. B-сплайн $B_{i,k}(x)$ степені k на вузловому векторі t_i визначається рекурсивно за формулою Кокса-де Бура [11]:

$$B_{i,0}(x) = \begin{cases} 1 & \text{якщо } t_i \leq x < t_{i+1} \\ 0 & \text{інакше} \end{cases} \quad (2)$$

$$B_{i,k}(x) = \frac{x - t_i}{t_{i+k} - t_i} B_{i,k-1}(x) + \frac{t_{i+k+1} - x}{t_{i+k+1} - t_{i+1}} B_{i+1,k-1}(x) \quad (3)$$

В-сплайни мають важливі властивості: локальну підтримку (зміна одного коефіцієнта впливає лише локально) [11], невід'ємність [11], розбиття одиниці $\sum_i B_{i,k}(x) = 1$ [11], та гладкість порядку $k - 1$ [11]. Локальна природа В-сплайнів забезпечує розрідженість обчислювальних графів, що є критичним для ефективності KAN [8].

2.3. Архітектура мереж Колмогорова-Арнольда

KAN-шар відображає вектор з n_{in} входів на вектор з n_{out} виходів через навчені одновимірні функції [8]:

$$y_j = \sum_{i=1}^{n_{in}} \varphi_{j,i}(x_i), \quad j = 1, \dots, n_{out}, \quad (4)$$

де кожна функція $\varphi_{j,i}$ параметризується як комбінація базисної активації та В-сплайну [8]:

$$\varphi_{j,i}(x) = w_b \cdot b(x) + w_s \cdot \sum_{l=0}^{G+k-1} c_{j,i,l} B_{l,k}(x), \quad (5)$$

де $b(x)$ – базова функція активації (зазвичай SiLU) [8], w_b, w_s – навчені скалярні ваги [8], $c_{j,i,l}$ – коефіцієнти сплайну [8], G – розмір сітки, k – степінь сплайну [8].

3. МЕТОД KAN-КЕРОВАНОГО ПРУНІНГУ

3.1. Гібридна архітектура CNN-KAN

Пропонується гібридна архітектура, що поєднує переваги згорткових мереж для просторової обробки та KAN для семантичного аналізу [8]. Архітектура складається з чотирьох компонентів:

Згортковий backbone генерує C карт ознак розміром " $H \times W$ " через послідовність residual блоків [1]. Global Average Pooling перетворює просторові карти у вектор $f \in R^C$:

$$f_c = \frac{1}{H \cdot W} \sum_{h=1}^H \sum_{w=1}^W A_c[h, w], \quad c = 1, \dots, C \quad (6)$$

KAN-bottleneck відображає вектор ознак у латентний простір меншої розмірності $L < C$ [8]:

$$z_j = \sum_{i=1}^c \varphi_{j,i}(f_i), \quad j = 1, \dots, L \quad (7)$$

Лінійний класифікатор відображає латентний вектор на логіти для K класів.

3.2. Критерій важливості фільтрів

Важливість i -го згорткового фільтра оцінюється через його вплив на латентний простір KAN. Для кожної пари (вхід i , латентний нейрон j) сила впливу визначається L2-нормою коефіцієнтів сплайну:

$$I_{ji} = |w_{s,ji}| \cdot \sqrt{\sum_{l=0}^{G+k-1} c_{ji,l}^2} \quad (8)$$

Загальна важливість i -го фільтра визначається як його максимальний вплив на будь-який латентний нейрон:

$$\text{Importance}(i) = \max_{j=1,\dots,L} I_{ji} \quad (9)$$

Використання максимуму обґрунтовується спеціалізацією латентних нейронів: якщо фільтр має сильний вплив хоча б на один нейрон, він кодує унікальну семантичну інформацію.

3.3. Алгоритм ітеративного прунінгу

Запропоновано алгоритм ітеративного структурного прунінгу:

1. Навчити повну CNN-KAN модель до збіжності
2. Для кожної ітерації прунінгу:
 - Обчислити важливість усіх фільтрів за формулами (8–9).
 - Видалити $p\%$ найменш важливих фільтрів.
 - Структурно модифікувати архітектуру.
 - Донавчити модель протягом E епох.
3. Повторювати до досягнення цільового рівня стиснення або падіння точності.

4. ЕКСПЕРИМЕНТАЛЬНІ РЕЗУЛЬТАТИ

4.1. Налаштування експериментів

Валідацію проведено на датасеті CIFAR-10 [12], що містить 50,000 навчальних та 10,000 тестових зображень 32×32 пікселі у 10 класах. Використано аугментацію (random crop, horizontal flip, color jitter) та нормалізацію за статистиками датасету. Параметри навчання: optimizer Adam, початкова швидкість 0.01, weight decay 5×10^{-4} , batch size 128, 100 епох з cosine annealing. Параметри KAN: grid size 5, spline order 3, латентна розмірність $L = C/4$.

4.2. Порівняння архітектур

Спочатку порівняно базові архітектури з KAN-модифікаціями (табл. 1).

Таблиця 1. Порівняння базових архітектур на CIFAR-10

Модель	Точність, %	Параметри	Час інференсу, мс
ResNet-18	93.8	11.2M	12.3
ResNet-18-KAN	94.1	11.5M	14.2
VGG-16-KAN	92.7	14.8M	18.5
WRN-28-10-KAN	95.2	36.5M	28.7

ResNet-18-KAN демонструє кращу точність (+0.3%) при незначному збільшенні параметрів, що підтверджує ефективність KAN-шарів.

4.3. Результати прунінгу

Застосовано ітеративний прунінг до ResNet-18-KAN (табл. 2).

Таблиця 2. Результати KAN-керованого прунінгу ResNet-18-KAN

Ітерація	Точність, %	Параметри	Compression	Speedup
0 (baseline)	94.1	11.5M	1.00x	1.00x
1 (5% pruned)	94.0	10.9M	1.05x	1.08x
2 (10% pruned)	93.9	10.4M	1.11x	1.15x
3 (15% pruned)	93.8	9.8M	1.17x	1.24x
4 (20% pruned)	93.5	9.2M	1.25x	1.32x
5 (25% pruned)	92.7	8.6M	1.34x	1.41x

При 25% прунінгу втрата точності становить лише 1.4%, що є прийнятним для 34% стиснення та 41% прискорення.

4.4. Порівняння з іншими методами

Проведено порівняння з класичними методами стиснення (табл. 3).

Таблиця 3. Порівняння методів стиснення на ResNet-18-KAN

Метод	Точність, %	Compression	Speedup	Параметри
Baseline	94.1	1.00x	1.00x	11.5M
Magnitude pruning	92.8	1.33x	1.38x	8.6M
Knowledge distillation	93.5	1.28x	1.35x	9.0M
KAN-guided pruning	94.2	1.31x	1.37x	8.8M

KAN-керований прунінг досягає найкращої точності (94.2%), перевершуючи magnitude-based підхід на 1.4% та distillation на 0.7% при порівнянному рівні стиснення.

5. ВИСНОВКИ

Запропоновано новий метод структурного прунінгу згорткових нейронних мереж, заснований на інтерпретабельності мереж Колмогорова-Арнольда. Ключовою інновацією є використання KAN-шару як інтерпретованого вузького місця для аналізу функціональної важливості згорткових фільтрів через оцінку їхнього впливу на латентний простір.

Розроблено математично обґрунтований критерій важливості на основі норми коефіцієнтів B-сплайнів та алгоритм ітеративного прунінгу з донавчанням. Метод дозволяє не лише видаляти параметри, але й пояснювати причини рішень через візуалізацію навчених функцій активації.

Експериментальна валідація на CIFAR-10 підтвердила ефективність методу. При 25% структурного прунінгу досягнуто 94.2% точності, що на 1.4% краще за magnitude-based pruning та на 0.7% краще за knowledge distillation. Метод забезпечує 34% зменшення розміру моделі та 41% прискорення інференсу при мінімальній втраті точності.

Перспективи подальших досліджень включають масштабування методу на більші датасети (ImageNet), застосування до сучасних архітектур (EfficientNet, Vision Transformers) та комбінування з іншими техніками стиснення.

ПЕРЕЛІК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. He K., Zhang X., Ren S., Sun J. Deep residual learning for image recognition. IEEE Conference on Computer Vision and Pattern Recognition, pp. 770-778, 2016.
2. Strubell E., Ganesh A., McCallum A. Energy and policy considerations for deep learning in NLP. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pp. 3645-3650, 2019.
3. Blalock D., Ortiz J.J.G., Frankle J., Gutttag J. What is the state of neural network pruning? Proceedings of Machine Learning and Systems, vol. 2, pp. 129-146, 2020.
4. Han S., Pool J., Tran J., Dally W. Learning both weights and connections for efficient neural networks. Advances in Neural Information Processing Systems, pp. 1135-1143, 2015.
5. Hinton G., Vinyals O., Dean J. Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531, 2015.
6. Jacob B., Kligys S., Chen B., et al. Quantization and training of neural networks for efficient integer-arithmetic-only inference. IEEE Conference on Computer Vision and Pattern Recognition, pp. 2704-2713, 2018.
7. Frankle J., Carbin M. The lottery ticket hypothesis: Finding sparse, trainable neural networks. International Conference on Learning Representations (ICLR), 2019.
8. Liu Z., Wang Y., Vaidya S., et al. KAN: Kolmogorov-Arnold Networks. arXiv preprint arXiv:2404.19756, 2024.
9. Kolmogorov A.N. On the representation of continuous functions of many variables by superposition of continuous functions of one variable and addition. Doklady Akademii Nauk SSSR, vol. 114, no. 5, pp. 953-956, 1957.
10. Arnold V.I. On functions of three variables. Doklady Akademii Nauk SSSR, vol. 114, pp. 679-681, 1957.
11. de Boor C. On calculating with B-splines. Journal of Approximation Theory, vol. 6, no. 1, pp. 50-62, 1972.
12. Krizhevsky A., Hinton G. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.