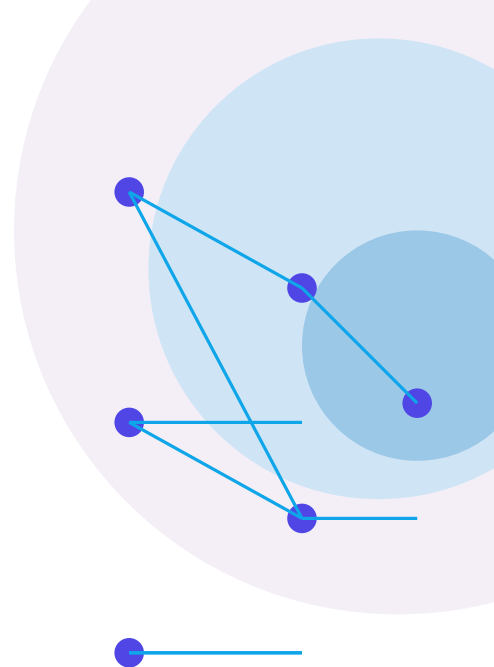


НЕЙРОННІ МЕРЕЖІ

Від нейрона до глибоких архітектур: MLP · CNN · RNN



01

НЕЙРОН І ПРОСТІ МЕРЕЖІ

- Біологічний та штучний нейрон
- Функції активації
- Перцептрон · MLP
- Backpropagation
- Функції втрат · Оптимізатори · Регуляризація

02

ЗГОРТКОВІ МЕРЕЖІ (CNN)

- Операція згортки та ядра
- Шари пулінгу
- ResNet, VGG, GoogLeNet
- Transfer Learning

03

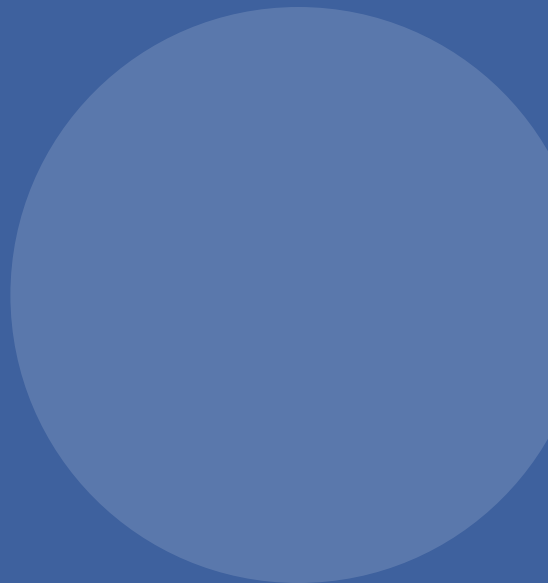
РЕКУРЕНТНІ МЕРЕЖІ (RNN)

- Класичний RNN · BPTT
- LSTM та GRU
- Seq2Seq · Attention
- Трансформери (огляд)

01

НЕЙРОН І ПРОСТІ МЕРЕЖІ

Від біологічного нейрона до багат шарового перцептрона



БІОЛОГІЧНИЙ НЕЙРОН

Дендрити

Отримують вхідні сигнали від інших нейронів. Аналог: входи x_1, x_2, \dots, x_n

Тіло клітини (сома)

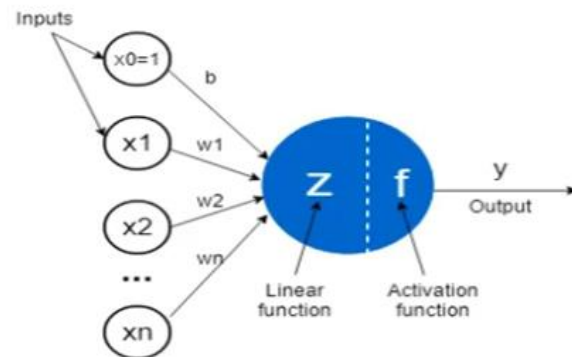
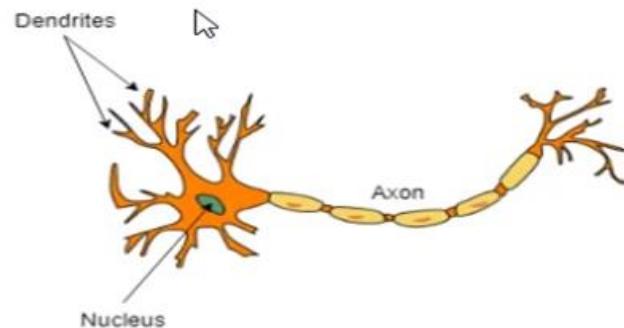
Підсумовує вхідні сигнали. Аналог: зважена сума $\sum w_i x_i + b$

Аксон

Передає вихідний сигнал іншим нейронам. Аналог: вихід $y = f(\text{net})$

Синапс

З'єднання між нейронами, що має певну «силу». Аналог: вага w_i



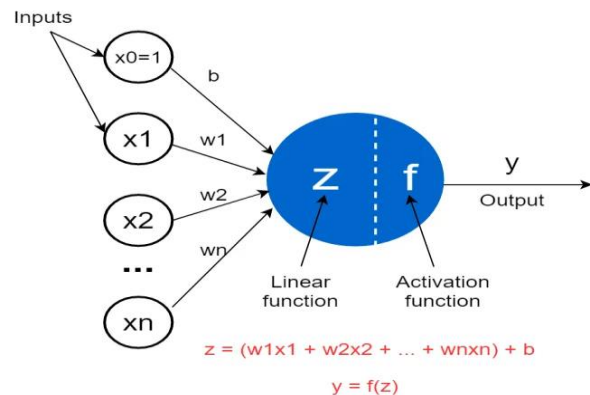
Штучний нейрон моделює основні функції біологічного: отримує сигнали, зважує їх і передає через функцію активації.

МАТЕМАТИЧНА МОДЕЛЬ НЕЙРОНА

$$\text{net} = \sum_i w_i \cdot x_i + b = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$y = f(\text{net}) \quad \text{де } f - \text{ функція активації}$$

$$\text{Матрична форма: } y = f(W^T x + b)$$



Позначення

x_i Вхідні ознаки (features)

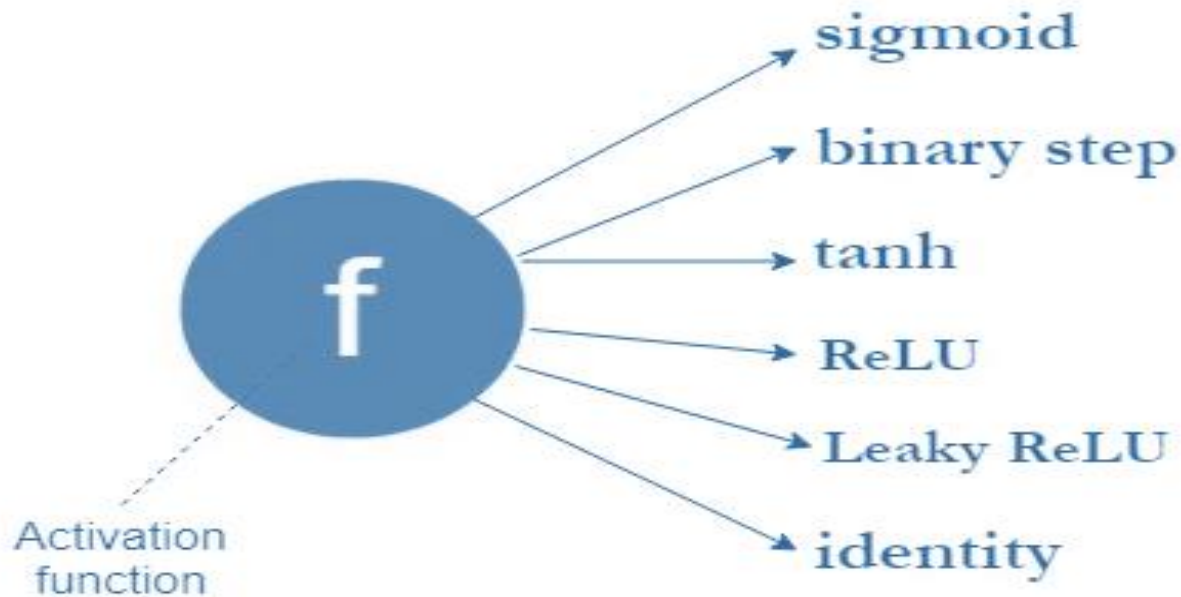
w_i Ваги — навчені параметри

b Зміщення (bias) — зсув порогу

net Зважена сума — лінійна комбінація

$f(\cdot)$ Функція активації — нелінійність

y Вихід нейрона



Sigmoid

$$\sigma(x) = 1 / (1 + e^{-x})$$

Діапазон: (0, 1)

+ Гладка, імовірнісна

- Vanishing gradient

Tanh

$$\tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$$

Діапазон: (-1, 1)

+ Центрована навколо 0

- Vanishing gradient

ReLU

$$\text{ReLU}(x) = \max(0, x)$$

Діапазон: [0, +∞)

+ Швидка, без насичення

- Dying ReLU проблема

Leaky ReLU

$$f(x) = \max(\alpha x, x), \quad \alpha \approx 0.01$$

Діапазон: (-∞, +∞)

+ Вирішує dying ReLU

- Гіперпараметр α

Softmax

$$\sigma_i = e^{x_i} / \sum_j e^{x_j}$$

Діапазон: (0,1), сума=1

+ Вихідний шар (класиф.)

- Лише фінальний шар

GELU

$$\text{GELU}(x) = x \cdot \Phi(x)$$

Діапазон: $\approx (-0.17, +\infty)$

+ Transformer default

- Обчислювально складна

Визначення

Найпростіший нейрон з бінарним виходом. Класифікує лінійно роздільні набори даних.

$$\hat{y} = 1 \quad \text{якщо} \quad \sum w_i x_i + b \geq \theta$$

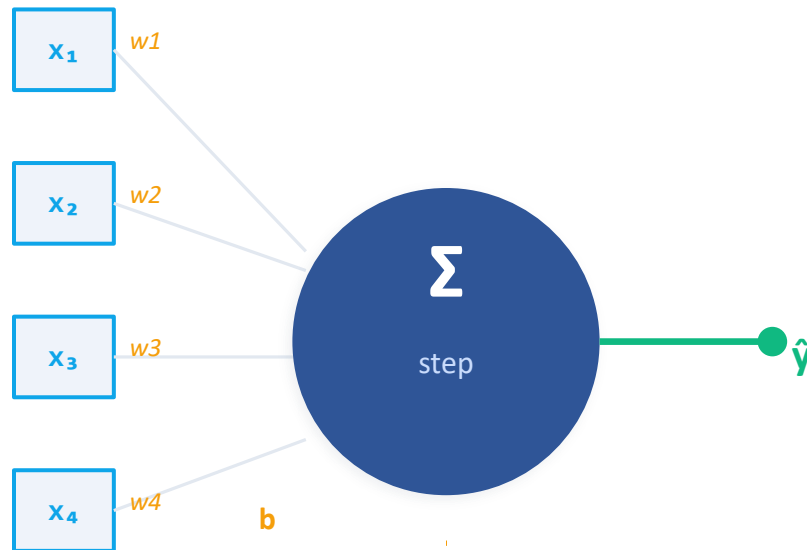
$$\hat{y} = 0 \quad \text{якщо} \quad \sum w_i x_i + b < \theta$$

Правило навчання (delta rule):

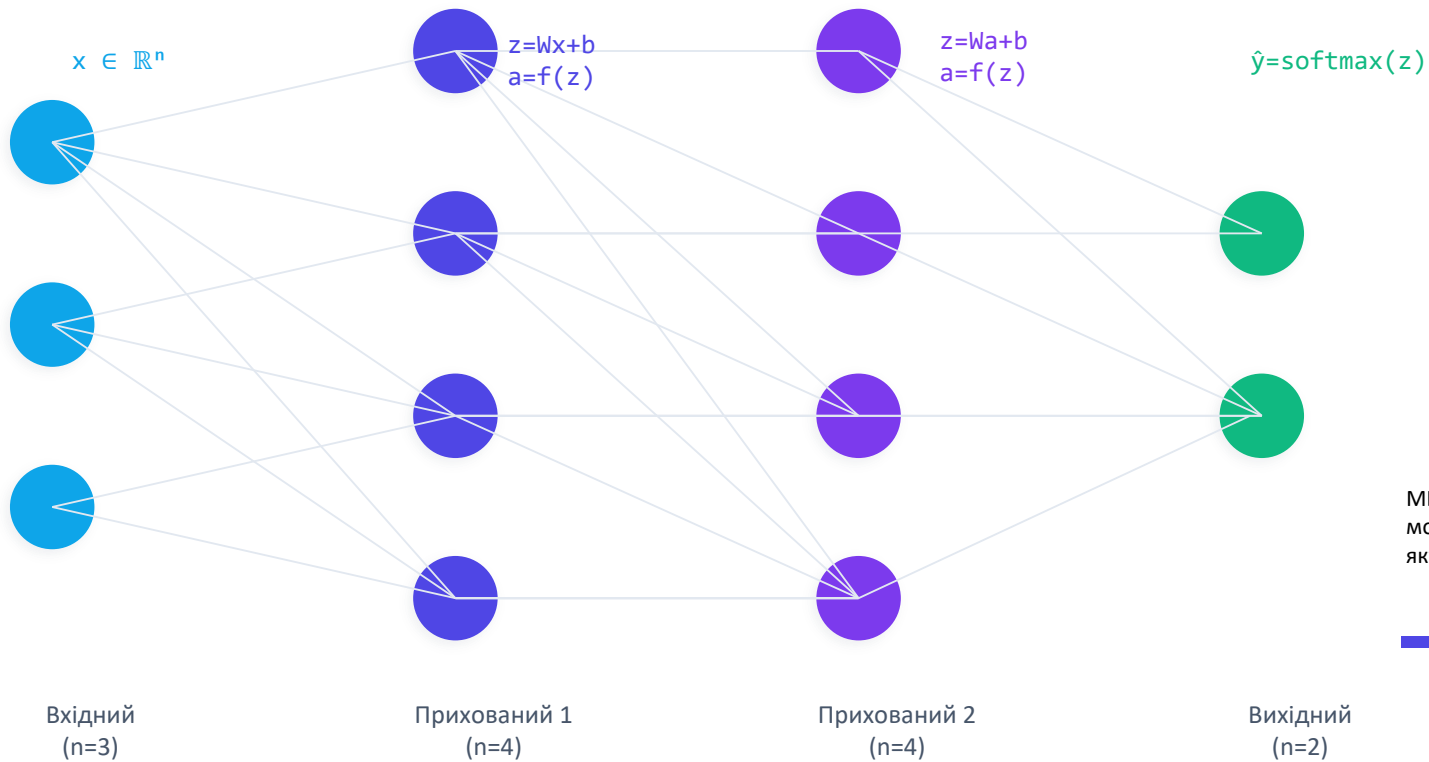
$$w_i \leftarrow w_i + \eta \cdot (y - \hat{y}) \cdot x_i$$

η — швидкість навчання (learning rate)

Обмеження: не може вирішити нелінійно нероздільні задачі (XOR).



БАГАТОШАРОВИЙ ПЕРЦЕПТРОН (MLP)



MLP з ≥ 1 прихованим шаром може апроксимувати будь-яку неперервну функцію.

ЗВОРОТНЕ ПОШИРЕННЯ ПОМИЛКИ (BACKPROPAGATION)

1

Forward Pass

Подаємо x , обчислюємо активації a^l по всіх шарах

$$a^l = f(W^l a^{l-1} + b^l)$$

2

Функція втрат

Порівнюємо вихід із правильною відповіддю

$$L = (1/m) \sum \text{loss}(\hat{y}_i, y_i)$$

3

Backward — вихідний шар δ^L

Похідна втрат по активаціях (delta-правило)

$$\delta^L = \nabla_a L \odot f'(z^L)$$

4

Backward — прихований шар δ^l

Поширюємо помилку назад через chain rule

$$\delta^l = ((W^{l+1})^T \delta^{l+1}) \odot f'(z^l)$$

5

Оновлення параметрів

Рухаємось у напрямку спадання градієнта

$$W^l \leftarrow W^l - \eta \cdot \delta^l (a^{l-1})^T$$

$$\text{Chain rule: } \partial L / \partial W^l = \delta^l \cdot (a^{l-1})^T$$

$$\partial L / \partial b^l = \delta^l$$

MSE

Mean Squared Error

$$(1/n)\sum(y_i - \hat{y}_i)^2$$

Задача: Регресія

Карає великі помилки

MAE

Mean Absolute Error

$$(1/n)\sum|y_i - \hat{y}_i|$$

Задача: Регресія

Стійка до викидів

**Binary
CE**

Binary Cross-Entropy

$$-(y \cdot \log \hat{y} + (1-y) \log(1-\hat{y}))$$

Задача: Бінарна кл.

Вихід: sigmoid

**Categ.
CE**

Categorical Cross-Entropy

$$-\sum_i y_i \cdot \log(\hat{y}_i)$$

Задача: Мульти-кл.

Вихід: softmax

Hinge

Hinge Loss

$$\sum \max(\theta, 1 - y_i \cdot \hat{y}_i)$$

Задача: SVM / класиф.

Великий відступ

KL Div.

Kullback–Leibler Divergence

$$\sum P(x) \cdot \log(P(x)/Q(x))$$

Задача: VAE · генеративні

Розбіжність розподілів

SGD

$$w \leftarrow w - \eta \cdot \nabla L$$

Стохастичний GD. Простий але нестабільний. Потребує ручного підбору η .

Momentum

$$\begin{aligned} v &\leftarrow \gamma v + \eta \cdot \nabla L \\ w &\leftarrow w - v \end{aligned}$$

Накопичує імпульс руху ($\gamma \approx 0.9$). Прискорює збіжність у стабільних напрямках.

RMSProp

$$\begin{aligned} E[g^2] &\leftarrow \rho E[g^2] + (1-\rho)g^2 \\ w &\leftarrow w - \eta / \sqrt{E[g^2] + \epsilon} \end{aligned}$$

Адаптивний η . Ефективний в RNN та онлайн-задачах.

Adam

$$\begin{aligned} m &\leftarrow \beta_1 m + (1-\beta_1)g \\ v &\leftarrow \beta_2 v + (1-\beta_2)g^2 \\ w &\leftarrow w - \eta \cdot \hat{m} / \sqrt{\hat{v} + \epsilon} \end{aligned}$$

Momentum + RMSProp. $\beta_1=0.9$, $\beta_2=0.999$. Найпопулярніший оптимізатор.

✓ Default

L2 (Weight Decay)

$$L_{total} = L + \lambda \cdot \sum w^2$$

Штрафує великі ваги. Запобігає overfitting.

L1 (Lasso)

$$L_{total} = L + \lambda \cdot \sum |w|$$

Обнуляє частину ваг → sparse мережа.

Dropout

$$p(\text{вимкнено}) = p \approx 0.3-0.5$$

Випадково вимикає нейрони під час навчання.

Early Stopping

Стоп якщо `val_loss` ↑ k епох

Зупиняємо при погіршенні val-метрики.

Batch Norm (BN)

$$\hat{x} = (x - \mu_B) / \sqrt{\sigma^2_B + \epsilon} \rightarrow \gamma \hat{x} + \beta$$

Нормалізує активації батчу. Прискорює навчання.

Layer Norm (LN)

$$\hat{x} = (x - \mu_L) / \sqrt{\sigma^2_L + \epsilon}$$

По ознаках (не батчу). Стандарт у Transformer.

У нейронній мережі прямого розповсюдження кожен вихід нейрону стає входом для наступного шару нейронів.

1) Вхід до прихованого шару обчислюється за формулою:

$$Z_{np} = \sum_{i=1}^N x_i w_i + b_1$$

де N - кількість нейронів у вхідному шарі (дорівнює кількості ознак)
 x_i - вхідні дані
 w_i - ваги
 b_1 - біас

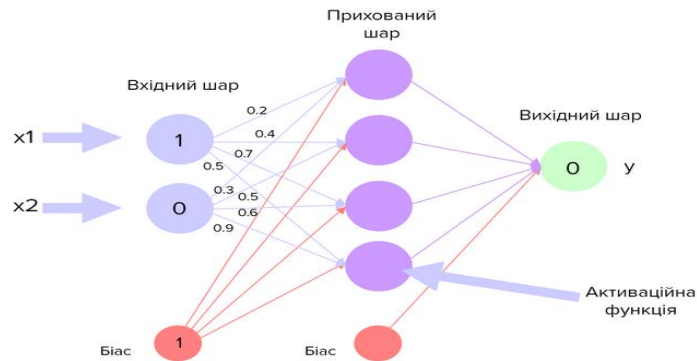
Нехай ваги та біас мають наступні значення:

$w_1=0.2$, $w_5=0.3$, $b_1=1$.

$w_2=0.4$, $w_6=0.5$,

$w_3=0.7$, $w_7=0.6$,

$w_4=0.5$, $w_8=0.9$,

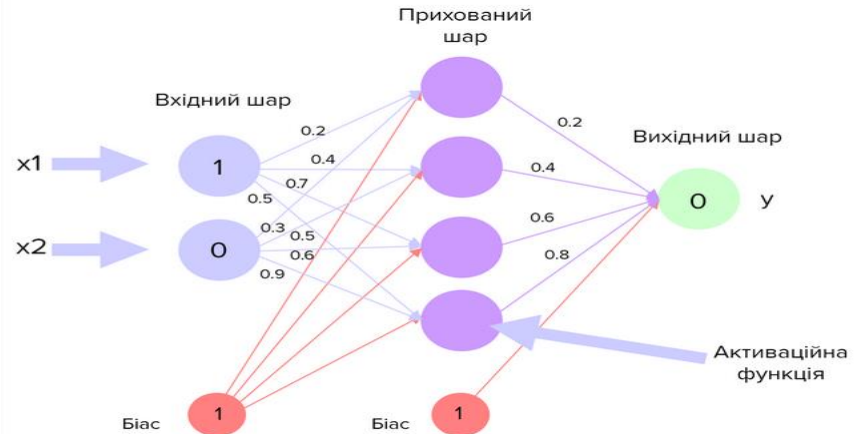
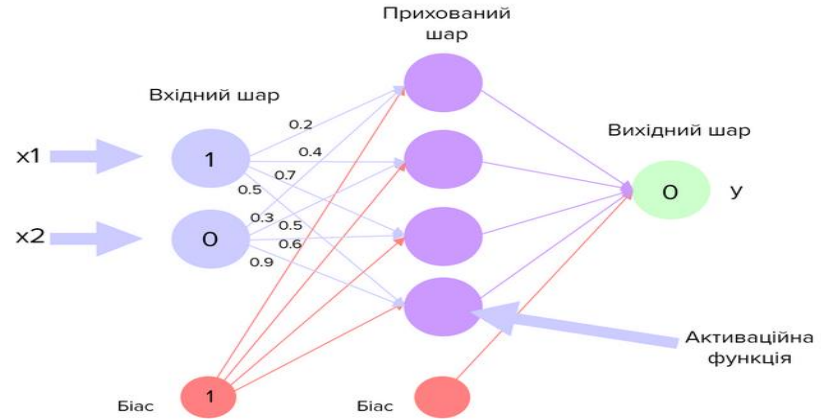


ПРИКЛАД

Кожен нейрон прихованого шару містить активаційну функцію. Вона може бути будь-якою, але будемо використовувати сігмоїд-функцію

Ця функція набуває значень на проміжку [0;1].

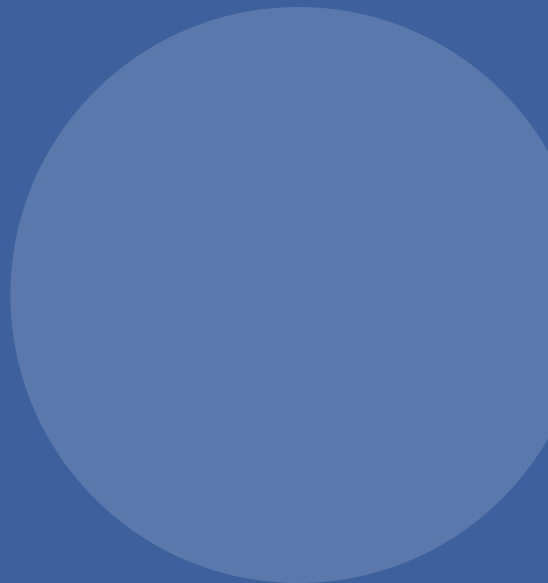
$$y_{np} = \frac{1}{1 + e^{-z_{np}}}$$



02

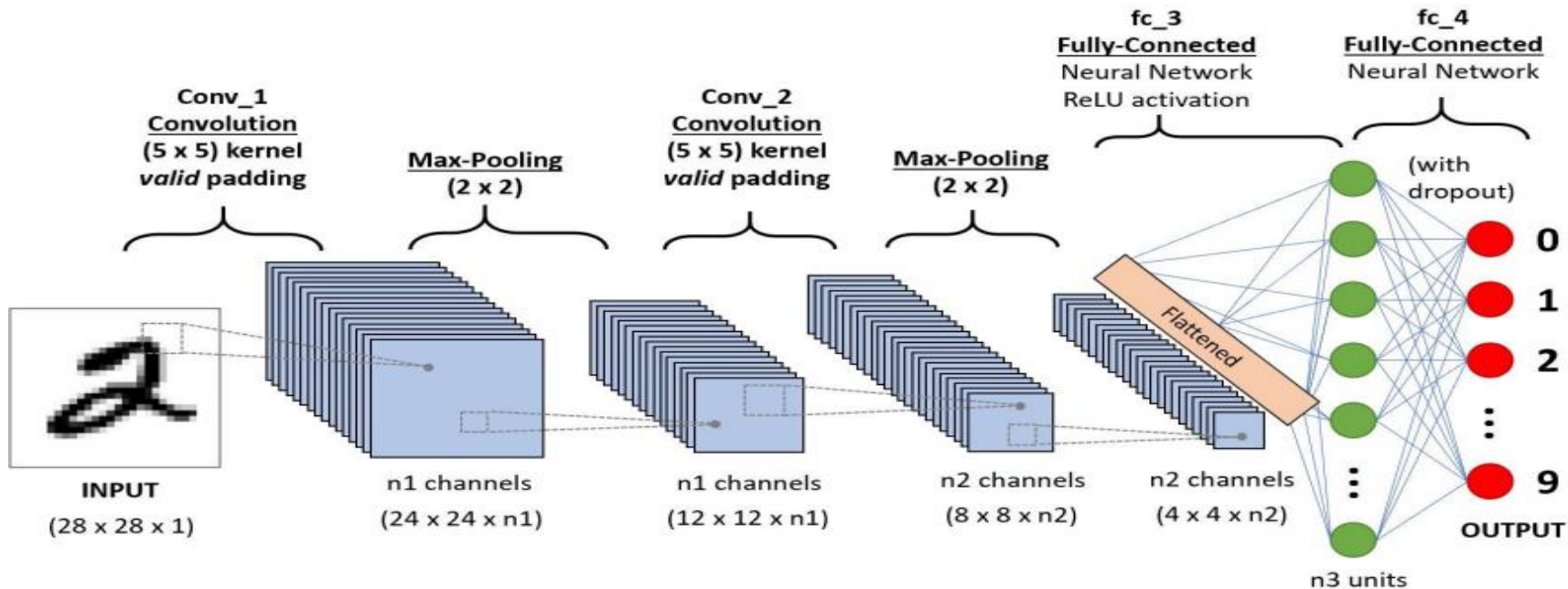
ЗГОРТКОВІ МЕРЕЖІ (CNN)

Просторова обробка даних: зображення, відео, сигнали



ЗГОРТКОВІ НЕЙРОННІ МЕРЕЖІ (CNN)

CNN — спеціалізована нейронна мережа для обробки даних з просторовою структурою (зображення, звук, текст). Використовує операцію згортки для локального видобутку ознак замість повно-з'єднаних шарів.



ЗАГАЛЬНИЙ ПОТІК: Вхід \rightarrow [Conv \rightarrow BN \rightarrow ReLU] \times N \rightarrow Pooling \times M \rightarrow Flatten \rightarrow FC \rightarrow Softmax \rightarrow Клас

Шар Згортки

Convolutional Layer

x ₀₀	x ₀₁	x ₀₂	x ₀₃
x ₁₀	x ₁₁	x ₁₂	x ₁₃
x ₂₀	x ₂₁	x ₂₂	x ₂₃
x ₃₀	x ₃₁	x ₃₂	x ₃₃

Вхід (4x4)

w ₀₀	w ₀₁	w ₀₂
w ₁₀	w ₁₁	w ₁₂
w ₂₀	w ₂₁	w ₂₂

Фільтр (3x3)



=

y ₀₀	y ₀₁
y ₁₀	y ₁₁

Карта ознак (2x2)

1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

Формула

$$(I * K)[i,j] = \sum \sum I[i+m,j+n,c] \cdot K[m,n,c] + b$$

Розмір виходу

$$W_{out} = \lfloor (W - F + 2P) / S \rfloor + 1$$

Параметри (3x3xCxK)

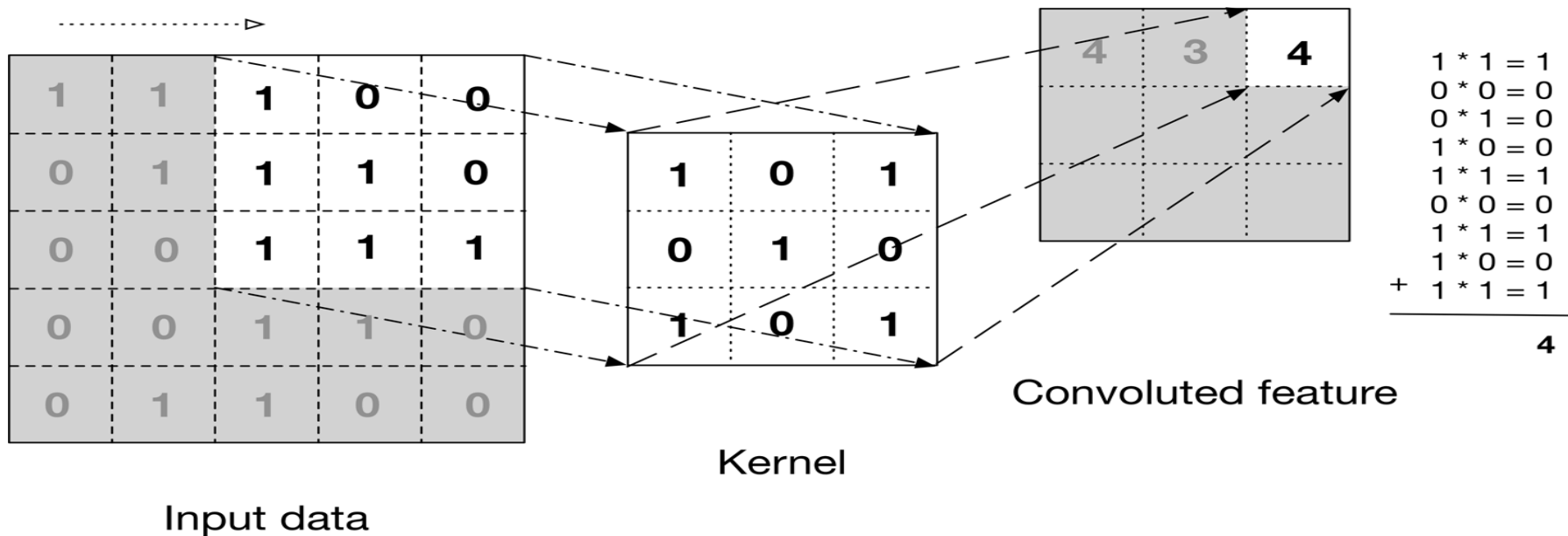
$$F \times F \times C_{in} \times K + K$$

Приклад (3x3, 32 фільтри, 3 канали)

$$3 \times 3 \times 3 \times 32 + 32 = 896 \text{ параметрів}$$

Шар Згортки

Convolutional Layer



Шар Згортки

Convolutional Layer

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

+

+ 1 = -25

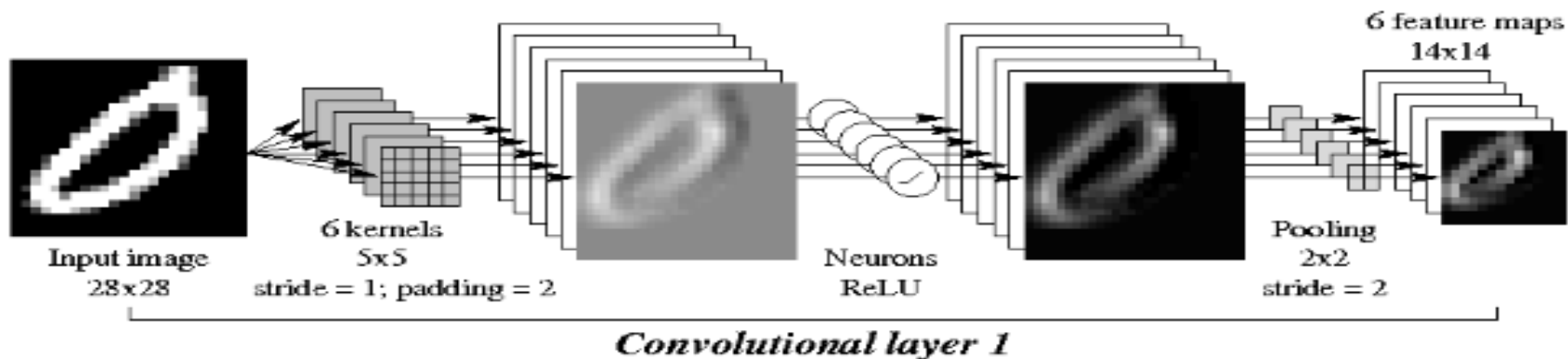
↑
Bias = 1

Output

-25				...
				...
				...
				...
...

Шар Згортки

Convolutional Layer



Згорткова нейронна мережа для задачі MNIST: детальний вигляд першого згорткового шару.

Max Pooling

$$\text{Out}[i,j] = \max(\text{вікно } 2 \times 2)$$

Вибирає максимум. Зберігає найяскравіші ознаки. Найпоширеніший.

Average Pooling

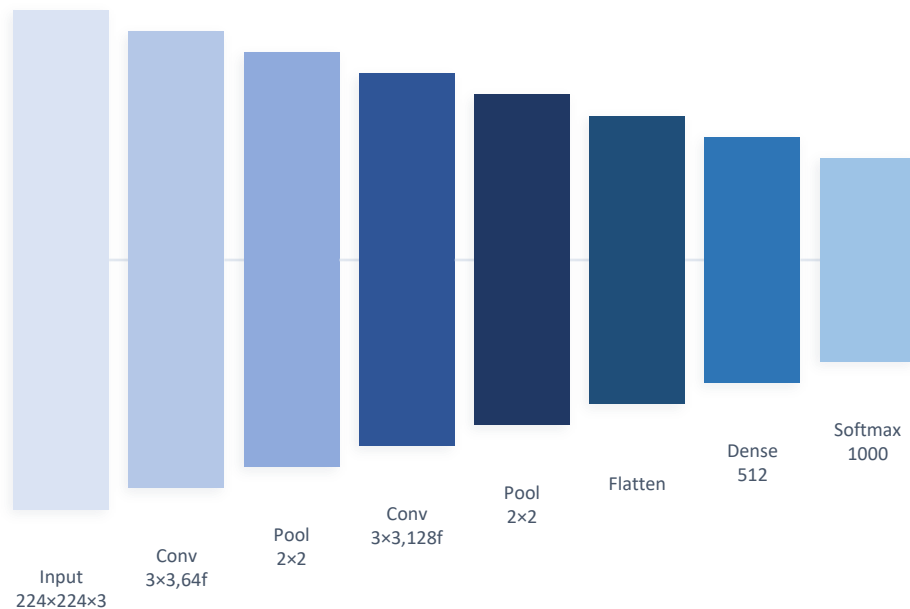
$$\text{Out}[i,j] = \text{mean}(\text{вікно } 2 \times 2)$$

Середнє вікно. М'якша версія, зберігає більше контексту.

Global Avg Pool

$$\text{Out} = \text{mean}(\text{карта } H \times W)$$

Замінює Dense-шар. Зменшує кількість параметрів у кінці CNN.



Принцип: Conv нарощують C_{out} (глибина), Pool зменшують $H \times W$ (простір). Dense робить класифікацію.

Convolution Output

1	2	2	1
4	9	1	0
1	5	2	3
4	6	1	2



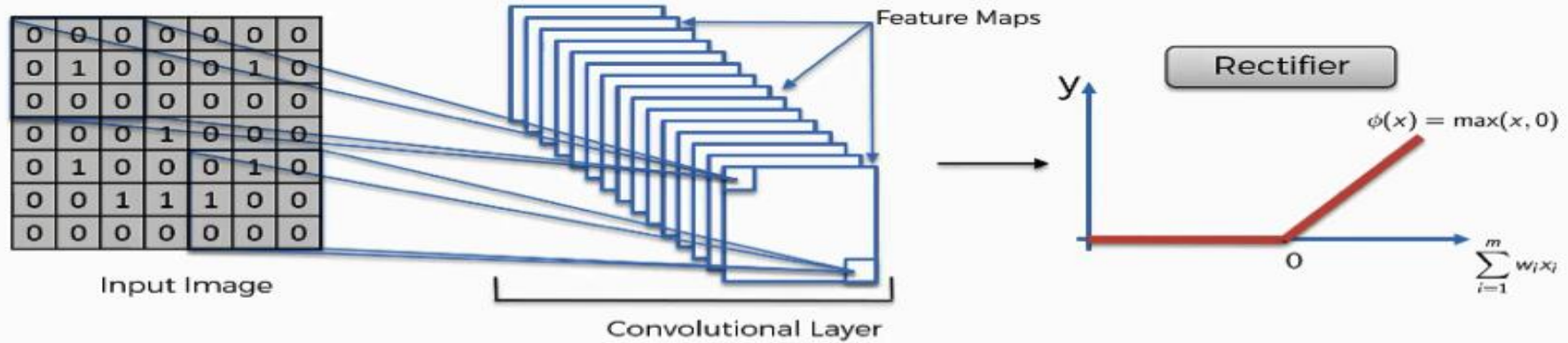
Max Pooling

9	2
6	3

Average Pooling

4	1
4	2





Математичне визначення функції активації ReLU -

$$f(x) = \max(0.0, x)$$

Dropout · Flatten · FC · Softmax

Dropout

Регуляризація

Під час навчання кожен нейрон деактивується з ймовірністю p ($0.2-0.5$). Запобігає перенавчанню.
 $f: y_i = x_i \cdot V_i / (1-p), \quad V_i \sim \text{Bernoulli}$

Flatten

Розгортання

3D тензор \rightarrow 1D вектор для FC-шару.
 $7 \times 7 \times 512 = 25,088$ значень
 ABO: GAP: $H \times W \times C \rightarrow C$ (512)

Fully Connected

Повноз'язний

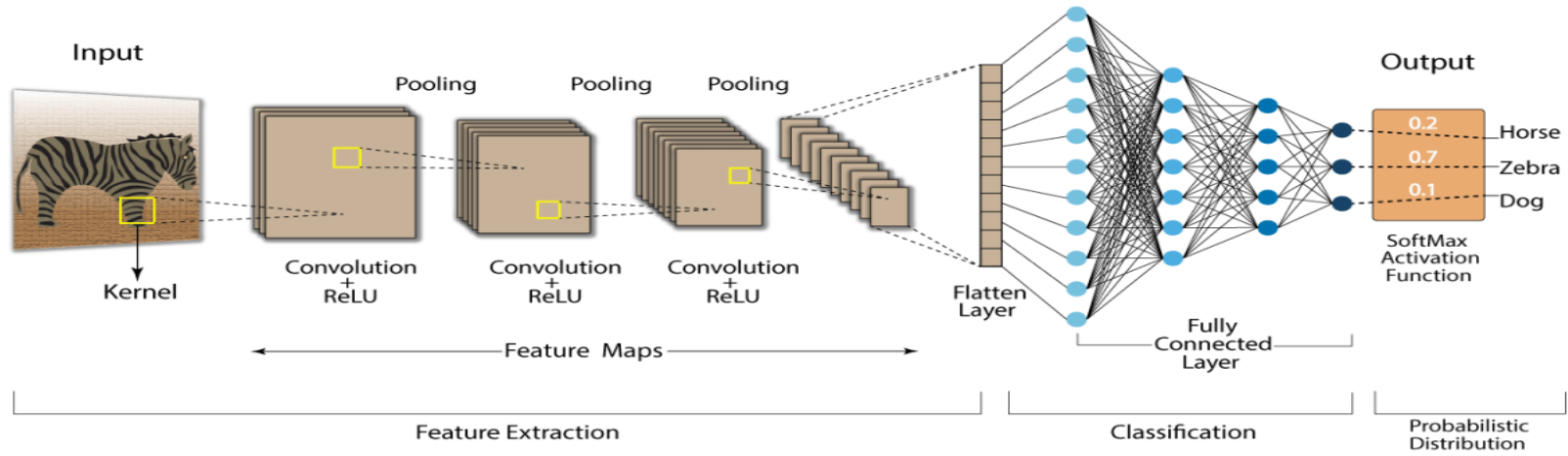
$y = f(Wx + b), \quad W \in \mathbb{R}^{(\text{out} \times \text{in})}$
 Параметри: $(\text{in}+1) \times \text{out}$
 $25088 \times 512 = 12,845,056$ ваг!

Softmax

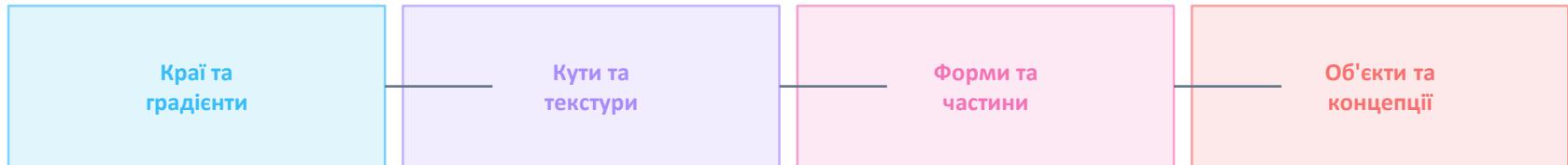
Вихідний шар

$\sigma(z)_i = e^{z_i} / \sum_j e^{z_j}$
 $\sum \sigma(z)_i = 1 \quad | \quad \text{Loss: } -\sum y_i \log(\hat{y}_i)$
 1000 класів (ImageNet)

Convolution Neural Network (CNN)



Ієрархія ознак:



1989

LeNet-5

Точність: ~99% MNIST

Парам: 60K

Перша успішна CNN (LeCun).
Conv+Pool+Dense.

2012

AlexNet

Точність: Top-5: 15.3%

Парам: 60M

Виграв ImageNet. ReLU, Dropout, GPU навчання.

2014

VGG-16

Точність: Top-5: 7.3%

Парам: 138M

Лише 3x3 фільтри, глибина 16–19 шарів.

2014

GoogLeNet

Точність: Top-5: 6.7%

Парам: 7M

Inception-модулі. Паралельні conv різних розмірів.

2015

ResNet-50

Точність: Top-5: 3.6%

Парам: 25M

Skip connections: $F(x)+x$. 152 шари.

2017

DenseNet

Точність: SOTA

Парам: 8M

Кожен шар з'єднаний з усіма попередніми.

2019

EfficientNet

Точність: Top-1: 84%+

Парам: 5–66M

Compound scaling: ширина + глибина + роздільн.

2021

Vision Transformer

Точність: 87%+

Парам: 86M+

Self-Attention замість Conv. Patches як токени.



Класифікація зображень

ImageNet, медична діагностика (МРТ, рентген), контроль якості на виробництві

SOTA: Top-1 ~90%+



Object Detection

YOLO, SSD, Faster R-CNN: знаходить і локалізує об'єкти з bounding boxes

SOTA: COCO mAP ~60%



Semantic Segmentation

U-Net, DeepLab: класифікація кожного пікселя. Медицина.

SOTA: mIoU ~85%



Face Recognition

FaceNet, ArcFace: ідентифікація та верифікація осіб (FaceID на телефонах)

SOTA: LFW ~99.8%



OCR / Document AI

Розпізнавання тексту з зображень, аналіз документів, перевірка підписів

SOTA: ~99% (друк)



Ігри та Reinforcement L.

Atari, AlphaGo — CNN обробляє стан поля, передає в policy network

SOTA: superhuman

Ідея: взяти модель, навчену на ImageNet (1.2М зображень), і донавчити на власних даних

1

Feature Extraction

Заморозити всі conv-шари. Замінити лише Dense-шар. Навчати тільки новий класифікатор.

Підходить: мало даних (100–1000 прикладів).
 $\eta = 0.001$

2

Fine-tuning (частковий)

Заморозити перші шари (загальні ознаки). Розморозити останні conv-блоки. Навчати з малим η .

Підходить: середня кількість даних.
 $\eta = 1e-5$

3

Full Fine-tuning

Розморозити всю мережу. Навчати з дуже малим learning rate. Потрібно багато даних.

Підходить: великий датасет (10K+).
 $\eta = 1e-6$

03

РЕКУРЕНТНІ МЕРЕЖІ (RNN)

Обробка послідовностей: текст, мовлення, часові ряди

Математична модель

RNN обробляє послідовності $x_t \in \mathbb{R}^n$ по одному кроку $t=1, \dots, T$:

$$h_t = \tanh(W_h \cdot h_{t-1} + W_x \cdot x_t + b)$$

$$\hat{y}_t = \text{softmax}(W_y \cdot h_t + b_y)$$

Позначення:

x_t — вхідний вектор на кроці t

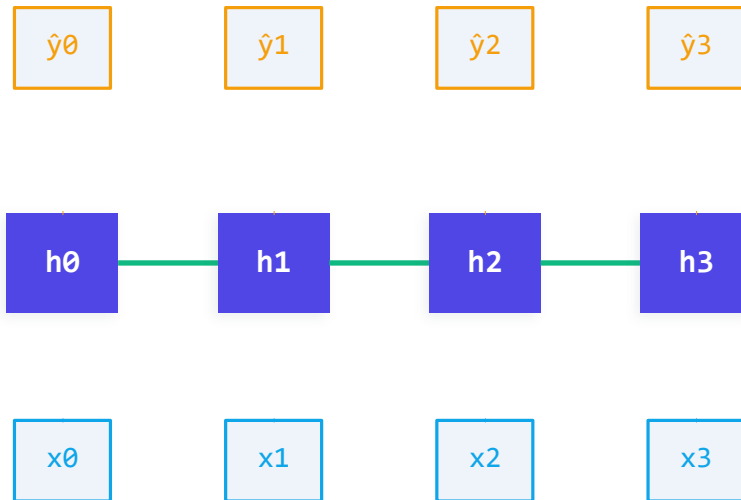
h_t — прихований стан — пам'ять

h_{t-1} — попередній прихований стан

W_h, W_x — спільні ваги для всіх t

\hat{y}_t — вихід на кроці t

Розгорнута RNN у часі (навч.: BPTT)



LONG SHORT-TERM MEMORY (LSTM, Hochreiter 1997)

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_t = o_t \odot \tanh(c_t)$$

f_t

Forget Gate

$$\sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Вирішує, яку частину c_{t-1} зберегти (0=забути, 1=залишити)

i_t

Input Gate

$$\sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

Скільки нової інформації записати до пам'яті

\tilde{c}_t

Cell Update

$$\tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

Новий кандидат у пам'ять (масштабується на i_t)

o_t

Output Gate

$$\sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

Скільки пам'яті видати як прихований стан h_t

Ключ LSTM: cell state c_t — «транспортна стрічка» — передає довгострокову пам'ять без перемноження градієнтів.

Рівняння GRU

Update Gate

$$z_t = \sigma(Wz \cdot [h_{t-1}, x_t])$$

Reset Gate

$$r_t = \sigma(Wr \cdot [h_{t-1}, x_t])$$

Candidate State

$$\tilde{h}_t = \tanh(W \cdot [r_t \odot h_{t-1}, x_t])$$

Output

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t$$

Update gate z_t : скільки попередн. стану зберегти. Reset gate r_t : скільки пам'яті враховувати.

LSTM vs GRU

	LSTM	GRU
Гейти	4	2
Стани	є	тільки h_t
Параметри	$\sim 4 \times H(H+D)$	$\sim 3 \times H(H+D)$
Пам'ять	Довша	Коротша
Швидкість	Повільніше	Швидше
Якість	Кращ. довгі	Кращ. короткі
Рік	1997	2014

$$\text{ВРТТ: } \partial L / \partial W_h = \sum_t \sum_{k \leq t} \left(\prod_{j=k+1}^t \partial h_j / \partial h_{j-1} \right) \cdot \partial h_k / \partial W_h$$

Vanishing Gradient

$$\|\partial h_t / \partial h_1\| = \|\prod_{\tau=1}^t \partial h_{\tau} / \partial h_{\tau-1}\| \rightarrow 0$$

Рішення

Добуток малих чисел (<1) \rightarrow градієнт зникає. Мережа не навчається на далеких залежностях.

Рішення: LSTM, GRU; Transformer; gradient clipping

Exploding Gradient

$$\|\partial h_t / \partial h_1\| \rightarrow \infty$$

Рішення: Gradient Clipping: $g \leftarrow g \cdot \max_norm / \|g\|$

Добуток великих значень \rightarrow NaN у градієнтах. Параметри розходяться під час навчання.

Послідовне навчання (складність)

$$O(T \cdot N^2) \text{ – немає паралелізму}$$

Рішення: Truncated BPTT; Transformer Architecture

RNN не можна навчати паралельно (залежність від попереднього кроку). Довго на довгих послідовн.



Машинний переклад

Seq2Seq LSTM перекладає речення (Google Translate до 2017). Many-to-Many.



Генерація тексту/музики

CharRNN передбачає наступний символ. Шекспір, стиль Баха.



Часові ряди

Прогноз акцій, погоди, енергоспоживання. Many-to-One або Many-to-Many.



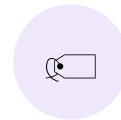
Аналіз тональності

Bi-LSTM: 'Мені сподобалось' → Positive. IMDB, Twitter.



Розпізнавання мовлення

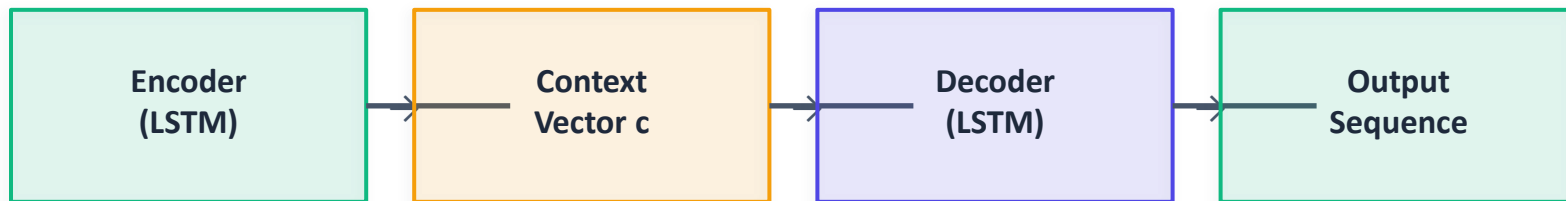
CTC+LSTM → акустична модель аудіо→текст (DeepSpeech, Whisper).



NER / POS tagging

Bi-LSTM-CRF: Named Entity Recognition — Особа, Організація, Місце.

SEQ2SEQ ТА МЕХАНІЗМ УВАГИ (ATTENTION)



Bahdanau Attention: $e_{ts} = \text{score}(s_t, h_s) \rightarrow \alpha_{ts} = \text{softmax}(e_{ts}) \rightarrow c_t = \sum_s \alpha_{ts} \cdot h_s$

Функції score (схожість запиту та ключа):

Dot-Product

$$\text{score}(s, h) = s^T h$$

Scaled Dot

$$\text{score} = s^T h / \sqrt{d_k}$$

Additive

$$\text{score} = v^T \cdot \tanh(W_s s + W_h h)$$

Cosine

$$\text{score} = s^T h / (\|s\| \cdot \|h\|)$$

Multi-Head: $MH(Q, K, V) = \text{Concat}(\text{head}_i) \cdot W^O$, $\text{head}_i = \text{Attention}(QW_i Q, KW_i K, VW_i V)$

Self-Attention: $A = \text{softmax}(QK^T / \sqrt{d_k}) \cdot V$ (Scaled Dot-Product)

TRANSFORMER: ATTENTION IS ALL YOU NEED (Vaswani 2017)

Self-Attention

Кожен токен зважає всі інші. Паралельне навчання без рекурентності.

$$A = \text{softmax}(QK^T/\sqrt{d_k})V$$

Positional Encoding

Додається до ембедінгів: sin/cos функції від позиції.

$$PE(\text{pos}, 2i) = \sin(\text{pos}/10000^{(2i/d)})$$

Feed-Forward Block

FFN застосовується до кожного токена незалежно: 2 шари + ReLU.

$$\text{FFN}(x) = \text{ReLU}(xW_1 + b_1)W_2 + b_2$$

Add & Layer Norm

Residual + нормалізація: стабілізує навчання глибоких стеків.

$$\text{Output} = \text{LayerNorm}(x + \text{Sublayer}(x))$$

Покоління Transformer

2017	Transformer	65M
2018	BERT	340M
2018	GPT-1	117M
2019	GPT-2	1.5B
2020	GPT-3	175B
2021	DALL-E	12B
2022	ChatGPT	~175B
2023	GPT-4/Gemini	?T

ПІДСУМКИ ЛЕКЦІЇ

Нейрон та MLP

● $net = \sum w_i x_i + b$, вихід $y = f(net)$

● Активації: ReLU (default), sigmoid, tanh, GELU

● MLP: $a^l = f(W^l a^{l-1} + b^l)$ — Universal Approximation Theorem

● Backprop + Chain Rule + Adam optimizer

CNN (Згорткові)

● Згортка: $(I * K)[i,j] = \sum \sum I[i+m, j+n] \cdot K[m,n]$

● Conv → BN → ReLU → Pool → Dense пайплайн

● ResNet: skip connections $F(x) + x$ — без vanish.grad

● Застосування: CV, медицина, OCR, сегментація

RNN / LSTM / GRU

● RNN: $h_t = \tanh(W_h h_{t-1} + W_x x_t + b)$ — послідовності

● LSTM: 4 гейти + cell state c_t — довга пам'ять

● GRU: 2 гейти (z, r) — ефективний аналог LSTM

● Transformer: Self-Attention замінив рекуренцію