

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN

CHUYÊN ĐỀ

AlexNet

MỤC LỤC

Chương 1	CNN Architectures	2
1.1	AlexNet Architecture	2
1.1.1	Kiến trúc mạng.....	2
1.1.2	INPUT.....	2
1.1.3	CONV1	2
1.1.4	NORM1.....	3
1.1.5	MAX POOL1	4
1.1.6	CONV2	4
1.1.7	NORM2.....	5
1.1.8	MAX POOL2.....	5
1.1.9	CONV3	6
1.1.10	CONV4	6
1.1.11	CONV5	7
1.1.12	MAX POOL3	8
1.1.13	FC6.....	8
1.1.14	FC7.....	8
1.1.15	FC8.....	9

CHƯƠNG 1 CNN ARCHITECTURES

1.1 AlexNet Architecture

1.1.1 Kiến trúc mạng

- Full (simplified) AlexNet architecture:
- [227x227x3] INPUT
- [55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
- [27x27x96] NORM1: Normalization layer
- [27x27x96] MAX POOL1: 3x3 filter at stride 2
- [27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
- [13x13x256] NORM2: Normalization layer
- [13x13x256] MAX POOL2: 3x3 filters at stride 2
- [13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1
- [13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1
- [13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
- [6x6x256] MAX POOL3: 3x3 filters at stride 2
- [4096] FC6: 4096 neurons
- [4096] FC7: 4096 neurons
- [1000] FC8: 1000 neurons (class scores)

1.1.2 INPUT

- INPUT: Ảnh đầu vào kích thước [227×227×3].

1.1.3 CONV1

- Input: Ảnh đầu vào M_1 kích thước [227×227×3].
- CONV1: gồm 96 filters với kích thước là [11×11×3] và bước nhảy là 4, pad là 0.
- Các giá trị thuộc filter này được tạo ngẫu nhiên.
- Output: Activation map M_2 có 96 feature map kích thước [55×55].
- Giải thích:

- + $\mathbf{W}_{\mathbf{M}_2} = \frac{(\mathbf{W}_{\mathbf{M}_1} - F + 2P)}{S} + 1 = \frac{(227 - 11 + 2 \cdot 0)}{4} + 1 = 55.$
 - $\mathbf{W}_{\mathbf{M}_1}$: Kích thước của ảnh đầu vào: 227.
 - F : Kích thước của filters: 11.
 - P : Padding: 0.
 - S : Bước nhảy của filters (Stride): 4.
 - $\mathbf{W}_{\mathbf{M}_2}$: Kích thước của activation map đầu ra: 55.
- ReLU activation function: $f(x) = \max(0, x)$
- $\mathbf{M}_2 = f(\mathbf{M}_2) = \max(0, \mathbf{M}_2)$
- Cụ thể:
 - + \mathbf{M}_2 : là activation map có 96 feature map với kích thước $[55 \times 55]$.
 - + Hàm ReLU được thực hiện trên tất cả phần tử nằm trong \mathbf{M}_2 .
 - + Giá trị trả về lại cho \mathbf{M}_2 .

1.1.4 NORM1

- Local Response Normalization (LRN):

$$b_{x,y}^i = \frac{a_{x,y}^i}{\left(k + \alpha \sum_{j=\max(0, i-\frac{n}{2})}^{\min(N-1, i+\frac{n}{2})} (a_{x,y}^j)^2 \right)^\beta}$$

- Giải thích:
 - + $a_{x,y}^i$: Sự hoạt động của 1 neuron tại filter i ở vị trí (x, y) .
 - + $b_{x,y}^i$: Sự hoạt động sau khi chuẩn hóa của 1 neuron tại kernel i ở vị trí (x, y) .
 - + Những siêu tham số cố định:
 - $k = 2$: Hằng số điều chỉnh độ lớn của mẫu chuẩn hóa. Có chức năng giúp tránh giá trị chuẩn hóa trở nên quá nhỏ giúp ổn định quá trình chuẩn hóa.
 - $n = 5$: Số filters lân cận của filter i .
 - $\alpha = 10^{-4}$: Hằng số điều chỉnh ảnh hưởng của các giá trị hoạt động đến mẫu chuẩn hóa. Giá trị α nhỏ hơn làm cho các hoạt động trước đó có ảnh hưởng lớn hơn đến kết quả chuẩn hóa, giúp kiểm soát mức độ "phạt" cho các giá trị lớn hơn.

- $\beta = 0.75$: Hằng số điều chỉnh độ dốc của hàm chuẩn hóa. Điều chỉnh cách thức mẫu chuẩn hóa ảnh hưởng đến các giá trị hoạt động.
- + Ý nghĩa: Sự hoạt động của 1 neuron tại filter i ở vị trí (x, y) được chia cho sự hoạt động của các neurons tại các filters lân cận ở vị trí (x, y) với các tỷ lệ của hằng số k, n, α, β .

1.1.5 MAX POOL1

- Input: Tensor \mathbf{M}_2 có 96 feature map kích thước $[55 \times 55]$.
- **MAX POOL1**: Sử dụng filter kích thước $[3 \times 3]$ với bước nhảy là 2 trên tất cả các filters trên activation map \mathbf{M}_2 .
- Output: Tensor \mathbf{M}_3 có 96 feature map kích thước $[27 \times 27]$.
- Giải thích:
 - + $\mathbf{W}_{\mathbf{M}_3} = \frac{(\mathbf{W}_{\mathbf{M}_2} - F + 2P)}{S} + 1 = \frac{(55 - 3 + 2 \cdot 0)}{2} + 1 = 27$.
 - $\mathbf{W}_{\mathbf{M}_2}$: Kích thước của tensor đầu vào: 55.
 - F : Kích thước của filters: 3.
 - P : Padding: 0.
 - S : Bước nhảy của filters (Stride): 2.
 - $\mathbf{W}_{\mathbf{M}_3}$: Kích thước của tensor đầu ra: 27.

1.1.6 CONV2

- Input: Tensor \mathbf{M}_3 có 96 feature map kích thước $[27 \times 27]$.
- **CONV2**: gồm 256 filters với kích thước $[5 \times 5 \times 96]$ và bước nhảy là 1, pad là 2.
- Các giá trị thuộc filter này được tạo ngẫu nhiên.
- Output: Activation map \mathbf{M}_4 có 256 feature map kích thước $[27 \times 27]$
- Giải thích:
 - + $\mathbf{W}_{\mathbf{M}_4} = \frac{(\mathbf{W}_{\mathbf{M}_3} - F + 2P)}{S} + 1 = \frac{(27 - 5 + 2 \cdot 2)}{1} + 1 = 27$.
 - $\mathbf{W}_{\mathbf{M}_3}$: Kích thước của tensor đầu vào: 27.
 - F : Kích thước của filters: 5.
 - P : Padding: 2.
 - S : Bước nhảy của filters (Stride): 1.
 - $\mathbf{W}_{\mathbf{M}_4}$: Kích thước của activation map đầu ra: 27.
 - ReLU activation function: $f(x) = \max(0, x)$
 - $\mathbf{M}_4 = f(\mathbf{M}_4) = \max(0, \mathbf{M}_4)$
 - Cụ thể:

- + \mathbf{M}_4 : là activation map có 256 feature map với kích thước $[27 \times 27]$.
- + Hàm ReLU được thực hiện trên tất cả phần tử nằm trong \mathbf{M}_4 .
- + Giá trị trả về lại cho \mathbf{M}_4 .

1.1.7 NORM2

- Local Response Normalization (LRN):

$$b_{x,y}^i = \frac{a_{x,y}^i}{\left(k + \alpha \sum_{j=\max(0, i-\frac{n}{2})}^{\min(N-1, i+\frac{n}{2})} (a_{x,y}^j)^2\right)^\beta}$$

- Giải thích:

- + $a_{x,y}^i$: Sự hoạt động của 1 neuron tại filter i ở vị trí (x, y) .
- + $b_{x,y}^i$: Sự hoạt động sau khi chuẩn hóa của 1 neuron tại kernel i ở vị trí (x, y) .
- + Những siêu tham số cố định:
 - $k = 2$: Hằng số điều chỉnh độ lớn của mẫu chuẩn hóa. Có chức năng giúp tránh giá trị chuẩn hóa trở nên quá nhỏ giúp ổn định quá trình chuẩn hóa.
 - $n = 5$: Số filters lân cận của filter i .
 - $\alpha = 10^{-4}$: Hằng số điều chỉnh ảnh hưởng của các giá trị hoạt động đến mẫu chuẩn hóa. Giá trị α nhỏ hơn làm cho các hoạt động trước đó có ảnh hưởng lớn hơn đến kết quả chuẩn hóa, giúp kiểm soát mức độ "phạt" cho các giá trị lớn hơn.
 - $\beta = 0.75$: Hằng số điều chỉnh độ dốc của hàm chuẩn hóa. Điều chỉnh cách thức mẫu chuẩn hóa ảnh hưởng đến các giá trị hoạt động.

Ý nghĩa: Sự hoạt động của 1 neuron tại filter i ở vị trí (x, y) được chia cho sự hoạt động của các neurons tại các filters lân cận ở vị trí (x, y) với các tỷ lệ của hằng số k, n, α, β .

1.1.8 MAX POOL2

- Input: Tensor \mathbf{M}_4 có 256 feature map kích thước $[27 \times 27]$.
- **MAX POOL2**: Sử dụng filter kích thước $[3 \times 3]$ với bước nhảy là 2 trên tất cả các filters trên activation map \mathbf{M}_4 .

- Output: Tensor \mathbf{M}_5 có 256 feature map kích thước $[13 \times 13]$.
- Giải thích:
 - + $\mathbf{W}_{\mathbf{M}_5} = \frac{(\mathbf{W}_{\mathbf{M}_4} - F + 2P)}{S} + 1 = \frac{(27 - 3 + 2 \cdot 0)}{2} + 1 = 13$.
 - $\mathbf{W}_{\mathbf{M}_4}$: Kích thước của tensor đầu vào: 27.
 - F : Kích thước của filters: 3.
 - P : Padding: 0.
 - S : Bước nhảy của filters (Stride): 2.
 - $\mathbf{W}_{\mathbf{M}_5}$: Kích thước của tensor đầu ra: 13.

1.1.9 CONV3

- Input: Tensor \mathbf{M}_5 có 256 feature map kích thước $[13 \times 13]$.
- **CONV3**: gồm 384 filters với kích thước là $[3 \times 3 \times 256]$.
- và bước nhảy là 1, pad là 1.
- Các giá trị thuộc filter này được tạo ngẫu nhiên.
- Output: Activation map \mathbf{M}_6 có 384 feature map kích thước $[13 \times 13]$.
- Giải thích:
 - + $\mathbf{W}_{\mathbf{M}_6} = \frac{(\mathbf{W}_{\mathbf{M}_5} - F + 2P)}{S} + 1 = \frac{(13 - 3 + 2 \cdot 1)}{1} + 1 = 13$.
 - $\mathbf{W}_{\mathbf{M}_5}$: Kích thước của tensor đầu vào: 13.
 - F : Kích thước của filters: 3.
 - P : Padding: 1.
 - S : Bước nhảy của filters (Stride): 1.
 - $\mathbf{W}_{\mathbf{M}_6}$: Kích thước của activation map đầu ra: 13.
- ReLU activation function: $f(x) = \max(0, x)$
- $\mathbf{M}_6 = f(\mathbf{M}_6) = \max(0, \mathbf{M}_6)$
- Cụ thể:
 - + \mathbf{M}_6 : là activation map có 384 feature map với kích thước $[13 \times 13]$.
- Hàm ReLU được thực hiện trên tất cả phần tử nằm trong \mathbf{M}_6 .
 - + Giá trị trả về lại cho \mathbf{M}_6 .

1.1.10 CONV4

- Input: Tensor \mathbf{M}_6 có 384 feature map kích thước $[13 \times 13]$.
- **CONV4**: gồm 384 filters với kích thước là $[3 \times 3 \times 384]$.
- và bước nhảy là 1, pad là 1.
- Các giá trị thuộc filter này được tạo ngẫu nhiên.

- Output: Activation map \mathbf{M}_7 có 384 feature map kích thước $[13 * 13]$.
 - Giải thích:
 - + $\mathbf{W}_{\mathbf{M}_7} = \frac{(\mathbf{W}_{\mathbf{M}_6} - F + 2P)}{S} + 1 = \frac{(13 - 3 + 2 \cdot 1)}{1} + 1 = 13$.
 - $\mathbf{W}_{\mathbf{M}_6}$: Kích thước của tensor đầu vào: 13.
 - F : Kích thước của filters: 3.
 - P : Padding: 1.
 - S : Bước nhảy của filters (Stride): 1.
 - $\mathbf{W}_{\mathbf{M}_7}$: Kích thước của activation map đầu ra: 13.
 - ReLU activation function: $f(x) = \max(0, x)$
 - $\mathbf{M}_7 = f(\mathbf{M}_7) = \max(0, \mathbf{M}_7)$
 - Cụ thể:
 - + \mathbf{M}_7 : là activation map có 384 feature map với kích thước $[13 \times 13]$.
 - Hàm ReLU được thực hiện trên tất cả phần tử nằm trong \mathbf{M}_7 .
- Giá trị trả về lại cho \mathbf{M}_7 .

1.1.11 CONV5

- Input: Tensor \mathbf{M}_7 có 384 feature map kích thước $[13 \times 13]$.
- **CONV5**: gồm 256 filters với kích thước là $[3 \times 3 \times 384]$.
- và bước nhảy là 1, pad là 1.
- Các giá trị thuộc filter này được tạo ngẫu nhiên.
- Output: Activation map \mathbf{M}_8 có 256 feature map kích thước $[13 \times 13]$.
- Giải thích:
 - + $\mathbf{W}_{\mathbf{M}_8} = \frac{(\mathbf{W}_{\mathbf{M}_7} - F + 2P)}{S} + 1 = \frac{(13 - 3 + 2 \cdot 1)}{1} + 1 = 13$.
 - $\mathbf{W}_{\mathbf{M}_7}$: Kích thước của tensor đầu vào: 13.
 - F : Kích thước của filters: 3.
 - P : Padding: 1.
 - S : Bước nhảy của filters (Stride): 1.
 - $\mathbf{W}_{\mathbf{M}_8}$: Kích thước của activation map đầu ra: 13.
 - ReLU activation function: $f(x) = \max(0, x)$
 - $\mathbf{M}_8 = f(\mathbf{M}_8) = \max(0, \mathbf{M}_8)$
 - Cụ thể:
 - + \mathbf{M}_8 : là activation map có 256 feature map với kích thước $[13 \times 13]$.
 - + Hàm ReLU được thực hiện trên tất cả phần tử nằm trong \mathbf{M}_8 .

+ Giá trị trả về lại cho \mathbf{M}_8 .

1.1.12 MAX POOL3

- Input: Tensor \mathbf{M}_8 có 256 feature map kích thước $[13 \times 13]$.
- **MAX POOL2**: gồm filter kích thước $[3 \times 3]$ với bước nhảy là 2 trên tất cả các filters trên activation map \mathbf{M}_8 .
- Output: Tensor \mathbf{M}_9 có 256 feature map kích thước $[6 \times 6]$.
- Giải thích:

$$+ \quad W_{\mathbf{M}_9} = \frac{(W_{\mathbf{M}_8} - F + 2P)}{S} + 1 = \frac{(13 - 3 + 2 \cdot 0)}{2} + 1 = 6.$$

- $W_{\mathbf{M}_9}$: Kích thước của tensor đầu vào: 13.
- F : Kích thước của filters: 3.
- P : Padding: 0.
- S : Bước nhảy của filters (Stride): 2.
- $W_{\mathbf{M}_8}$: Kích thước của tensor đầu ra: 6.

1.1.13 FC6

- Input: $H_1 = \mathbf{M}_9.\text{flatten}()$
- Output: Hidden layer \mathbf{H}_2 kích thước $[4096,]$.
- Giải thích:
 - + $\mathbf{H}_2 = W * H_1$.
 - \mathbf{H}_1 : Kích thước của activation map đầu vào sau khi duỗi
 - $W_{4096 \times 9216}$:
 - \mathbf{H}_2 : Kích thước của hidden layer đầu ra: $[4096,]$.
 - ReLU activation function: $f(x) = \max(0, x)$
 - $\mathbf{H}_2 = f(\mathbf{H}_2) = \max(0, \mathbf{H}_2)$
 - Cụ thể:
 - + \mathbf{H}_2 : là hidden layer có 4096 neuron.
 - + Hàm ReLU được thực hiện trên tất cả phần tử nằm trong \mathbf{H}_2 .
 - + Giá trị trả về lại cho \mathbf{H}_2 .
 - Dropout:
 - + Gán ngẫu nhiên 1 nửa số neuron của \mathbf{H}_2 với giá trị 0.

1.1.14 FC7

- Input: Hidden layer \mathbf{H}_2 kích thước $[4096,]$.
- Output: Hidden layer \mathbf{H}_3 kích thước $[4096,]$.
- Giải thích:

- + $\mathbf{H}_3 = W * \mathbf{H}_2$.
 - \mathbf{H}_2 : Kích thước của Hidden layer đầu vào
 - $W_{4096*4096}$.
 - \mathbf{H}_2 : Kích thước của hidden layer đầu ra: [4096,] .
- ReLU activation function: $f(x) = \max(0, x)$
- $\mathbf{H}_3 = f(\mathbf{H}_2) = \max(0, \mathbf{H}_2)$
- Cụ thể:
 - + \mathbf{H}_3 : là hidden layer có 4096 neuron.
 - + Hàm ReLU được thực hiện trên tất cả phần tử nằm trong \mathbf{H}_3 .
 - + Giá trị trả về lại cho \mathbf{H}_3 .
- Dropout:
 - + Gán ngẫu nhiên 1 nửa số neuron của \mathbf{H}_3 với giá trị 0.

1.1.15 FC8

- Input: Hidden layer \mathbf{H}_3 kích thước [4096,].
- Output: Hidden layer \mathbf{H}_4 kích thước [1000,].
- Giải thích:
 - + $\mathbf{H}_4 = W * \mathbf{H}_3$.
 - \mathbf{H}_3 : Kích thước của Hidden layer đầu vào
 - $W_{1000*4096}$.
 - \mathbf{H}_4 : Kích thước của hidden layer đầu ra: [1000,] .
- Softmax: $z(x) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$
- $\mathbf{H}_4 = z(\mathbf{H}_4) = \frac{e^{H_4^i}}{\sum_{j=1}^{1000} e^{H_4^j}}$
- Cụ thể:
 - + \mathbf{H}_4 : là hidden layer có 1000 neuron.
 - + Hàm Softmax được thực hiện trên tất cả phần tử nằm trong \mathbf{H}_4 .
 - + Giá trị trả về lại cho \mathbf{H}_4 .