

Basic CNN - Tutorial

Ngày 11 tháng 12 năm 2022

Phần I: Tổng quan Convolutional Neural Network (CNN)

1. Normal Neural Network trên images

Mỗi image có hàng ngàn pixels, mỗi pixel được xem như 1 feature, vậy nếu như một bức ảnh có kích thước $1000 * 1000$, thì sẽ có 1.000.000 features. Với normal feed-forward neural networks, mỗi layer là full-connected với previous input layer.

Trong normal feed-forward neural network, với mỗi layer, có 1.000.000 pixels, mỗi pixel lại kết nối full-connected với 1.000.000 pixels ở layer trước, tức sẽ có 1012 tham số. Đây là con số quá lớn để có thể tính được vào thời điểm đó, bởi vì với mô hình có nhiều tham số, sẽ dễ bị overfitted và cần lượng lớn data cho việc training, ngoài ra còn cần nhiều memory và năng lực tính toán cho việc training và prediction.

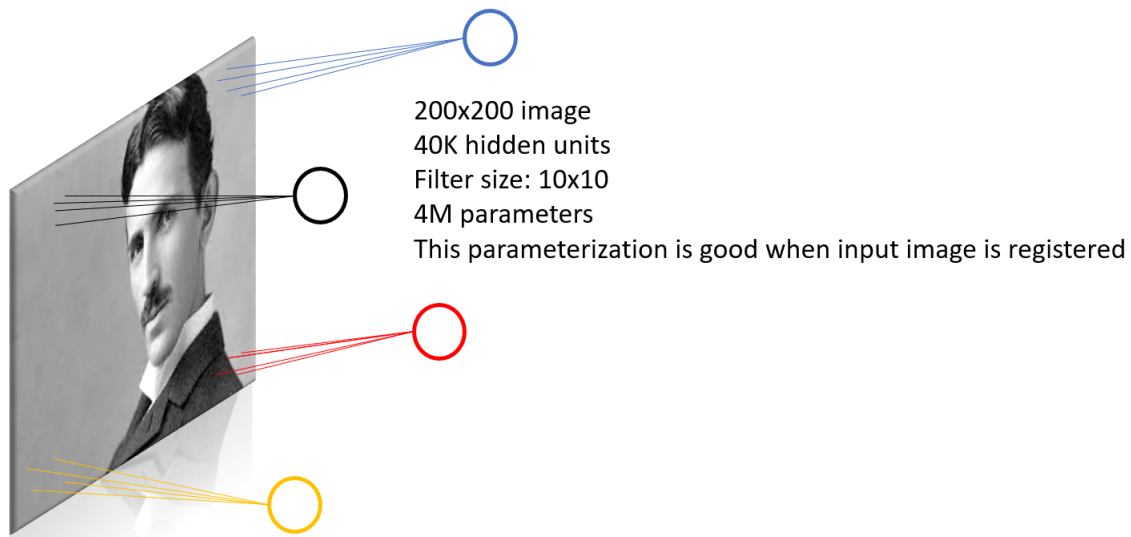
2. CNN hoạt động ra sao?

Có 2 đặc tính của image hình thành nên cách hoạt động của CNN trên image, đó là **feature localization** và **feature independence of location**.

Feature localization: mỗi pixel hoặc feature có liên quan với các pixel quanh nó.

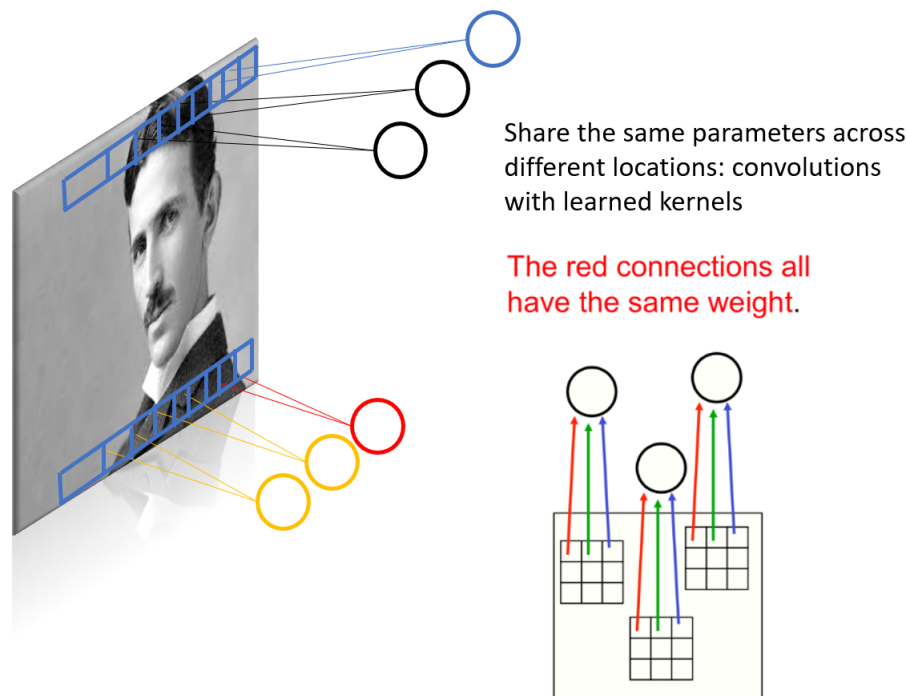
Feature Independence of location: mỗi feature dù nó có nằm ở đâu trong bức ảnh, thì nó vẫn mang giá trị của feature đó. CNN xử lý vấn đề có quá nhiều tham số với Shared parameters (feature independence of location) của Locally connected networks (feature localization), được gọi là Convolution Net.

Locally connected layer: Trong hidden layer đầu tiên, mỗi node sẽ kết nối tới một cụm nhỏ pixels của input image chứ không phải toàn bộ image, gọi là small portion. Theo cách này, ta sẽ có ít kết nối hơn, vì thế ít tham số hơn giữa input và hidden layer đầu tiên.



Hình 1: Local connected layer

Shared parameters: Có những khu vực, mà việc tìm ra feature là giống nhau về cách làm, vì vậy ta có thể dùng chung bộ parameter, trong hình trên là phía phải bên trên và phía trái bên dưới. Tức ta chia sẻ bộ parameter giữa những vị trí khác nhau trong bức ảnh.



Hình 2: Shared parameters

3. **Thành phần chính của CNN** Có 3 thành phần chính của CNN: convolution, pooling và fully-connected layer.

Convolution layer: gồm filter và Stride. Filter là locally connected network, mỗi filter (kernel) sẽ học được nhiều feature trong images. Mỗi filter sẽ di chuyển quanh bức ảnh với bước nhảy được cấu hình trước là Stride.

Pooling layer: bao gồm max-pooling và average-pooling layer tương ứng. Max-pooling layer chọn giá trị lớn nhất từ cửa sổ tính toán, còn Average-pooling layer sẽ tính giá trị trung bình của mỗi window.

Fully-connected layer: flatten convolution layer cuối cùng và fully connect neuron cho output layer.

4. Dataset - MNIST

Để đơn giản, chúng ta sẽ dùng bộ dataset MNIST, về phân lớp chữ biết tay, yêu cầu là cho một image, phân loại nó thuộc lớp nào trong 10 ký tự số.



Hình 3: Dataset MNIST

Mỗi image trong MNIST dataset có 28x28 ký tự dạng grayscale, được phân bố tại trung tâm bức ảnh.

Với một neural network thông thường, đã có thể xử lý tốt bài toán này. Bằng cách chuyển image 28x28 thành vector có $28 \times 28 = 784$ chiều, và đưa vector này vào input layer, đi qua một vài hidden layers, và kết thúc bởi output layer với 10 nodes, mỗi node đại diện cho 1 ký tự số.

Lưu ý rằng, với MNIST dataset chứa các image nhỏ, và ký tự số nằm ở trung tâm bức ảnh, chứ nếu ký tự nằm ở các góc cạnh của bức ảnh, sẽ gặp khó khăn. Điều này là khó khăn với những bức ảnh trong thực tế khi cần phân loại.

5. Các khái niệm quan trọng

(a) Convolutions

Convolutional neural networks cơ bản chỉ là convolutional layers, gọi tắt là conv layers, hoạt động dựa trên phép toán convolution của toán học. Các conv layers bao gồm các filters (hay còn gọi là convolutional kernel, convolutional filter, hoặc kernel), nó đơn giản là matrix 2 chiều gồm các con số. Ví dụ 3x3 filter.

Chúng ta có thể dùng 1 input image và 1 filter để tính output image bởi tích chập (convolving) filter với input image. Bao gồm các bước sau:

- Đặt filter lên phía trên image tại vài vị trí.
- Thực hiện element-wise multiplication giữa các giá trị trong filter và giá trị tương ứng trên image.
- Cộng các element-wise products lại. Kết quả tổng sẽ là output value cho destination pixel trên output image.
- Lặp lại cho các vị trí khác.

-1	0	1
-2	0	2
-1	0	1

Hình 4: 3x3 filter

Có thể mô phỏng 4 bước trên bằng một ví dụ nhỏ, với bức ảnh 4x4 grayscale và 3x3 filter:
 Các con số ở trên biểu diễn giá trị của pixel, trong khoảng từ 0 (black) đến 255 (white). Kết

0	50	0	29
0	80	31	2
33	90	0	75
0	9	0	95

-1	0	1
-2	0	2
-1	0	1

Hình 5: 4x4image-3x3filter

quả của phép tính input image và filter là matrix 2x2 có dạng như sau.

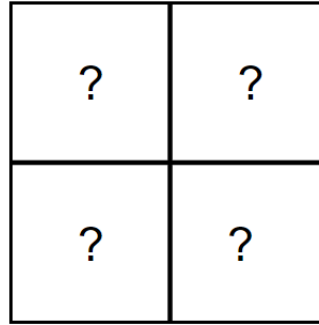
Đầu tiên ta đặt filter lên góc trái phía trên của input image:

Kết quả thu được là:

Tiếp theo cộng các kết quả lại: $62 - 33 = 29$, và đặt nó vào output image.

Trường hợp là ảnh màu, thì input là tensor 3 chiều, vì vậy filter là tensor 3 chiều trượt hết trong khối input.

Nếu ta muốn detect multiple features, ví dụ như muốn tìm cả horizontal edges và vertical



Hình 6: cnn-2x2output-image

0	50	0	29	-1	0	1
0	80	31	2	-2	0	2
33	90	0	75	-1	0	1
0	9	0	95			

Hình 7: cnn-input-output

Image Value	Filter Value	Result
0	-1	0
50	0	0
0	1	0
0	-2	0
80	0	0
31	2	62
33	-1	-33
90	0	0
0	1	0

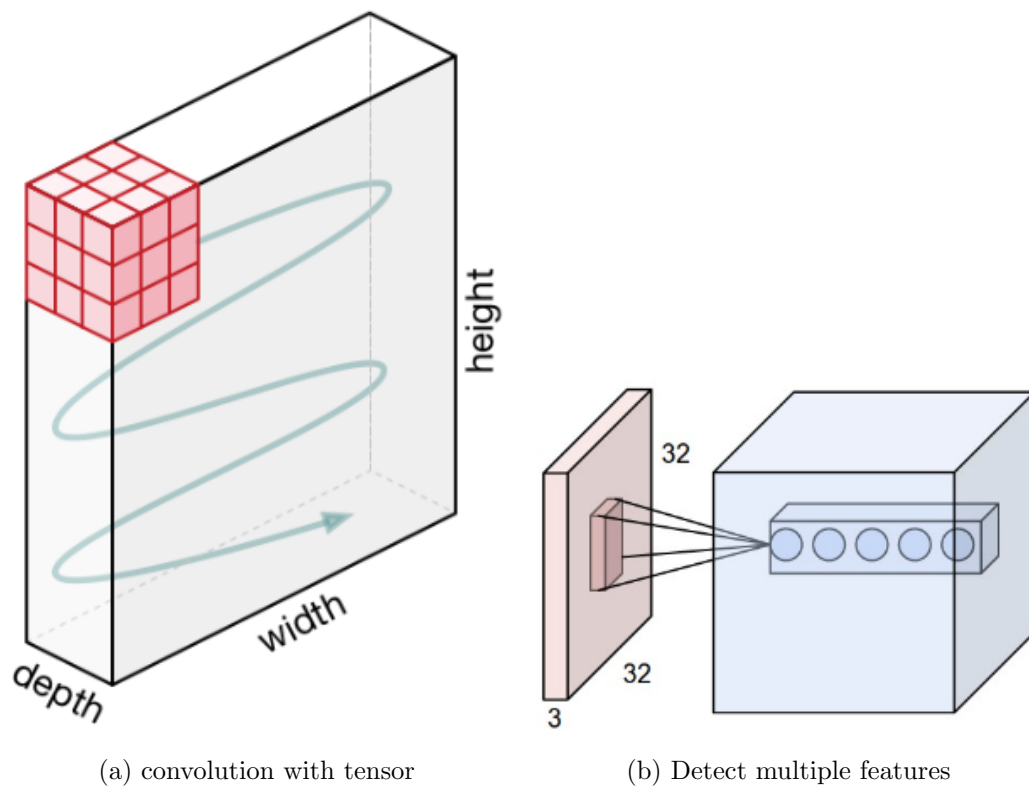
Hình 8: Bảng tính giá trị cho 1 pixel trên output image

edges, hoặc bất kỳ feature nào, ta sẽ muốn học từ nhiều convolutional filters trên cùng input. Đó là ta muốn tính toán đồng thời cho các kernel khác nhau.

0	50	0	29
0	80	31	2
33	90	0	75
0	9	0	95

29	-192
?	?

Hình 9: Output image



(a) convolution with tensor

(b) Detect multiple features

Hình 10: Convolution with tensor and multiple feature

```

1 class Conv3x3:
2     # A Convolution layer using 3x3 filters.
3
4     def __init__(self, num_filters):
5         self.num_filters = num_filters
6

```

```

7     # filters is a 3d array with dimensions (num_filters, 3, 3)
8     # We divide by 9 to reduce the variance of our initial values
9     self.filters = np.random.randn(num_filters, 3, 3) / 9
10
11 def iterate_regions(self, image):
12     # Generates all possible 3x3 image regions using valid padding. Image is
13     # a 2d numpy array.
14
15     h, w = image.shape
16
17     for i in range(h - 2):
18         for j in range(w - 2):
19             im_region = image[i:(i + 3), j:(j + 3)]
20             yield im_region, i, j
21
22 def forward(self, input):
23     # Performs a forward pass of the conv layer using the given input.
24     # Returns a 3d numpy array with dimensions (h, w, num_filters). Input is
25     # a 2d numpy array
26
27     self.last_input = input
28
29     h, w = input.shape
30     output = np.zeros((h - 2, w - 2, self.num_filters))
31
32     for im_region, i, j in self.iterate_regions(input):
33         output[i, j] = np.sum(im_region * self.filters, axis=(1, 2))
34
35     return output
36
37 def backprop(self, d_L_d_out, learn_rate):
38     # Performs a backward pass of the conv layer. d_L_d_out is the loss
39     # gradient for this layer's outputs, learn_rate is a float.
40
41     d_L_d_filters = np.zeros(self.filters.shape)
42
43     for im_region, i, j in self.iterate_regions(self.last_input):
44         for f in range(self.num_filters):
45             d_L_d_filters[f] += d_L_d_out[i, j, f] * im_region
46
47     # Update filters
48     self.filters -= learn_rate * d_L_d_filters
49
50     # We aren't returning anything here since we use Conv3x3 as the first
51     # layer in our CNN.
52     # Otherwise, we'd need to return the loss gradient for this layer's
53     # inputs, just like every
54     # other layer in our CNN.
55     return None

```

(b) **Filter in computer vision**

Để ý rằng, filter 3x3 ở trên, trong thị giác máy tính, gọi là vertical Sobel filter, và ta có thêm một dạng khác là horizontal Sobel filter.

Khi áp dụng mỗi loại filter trên vào input image, output image sẽ có dạng:

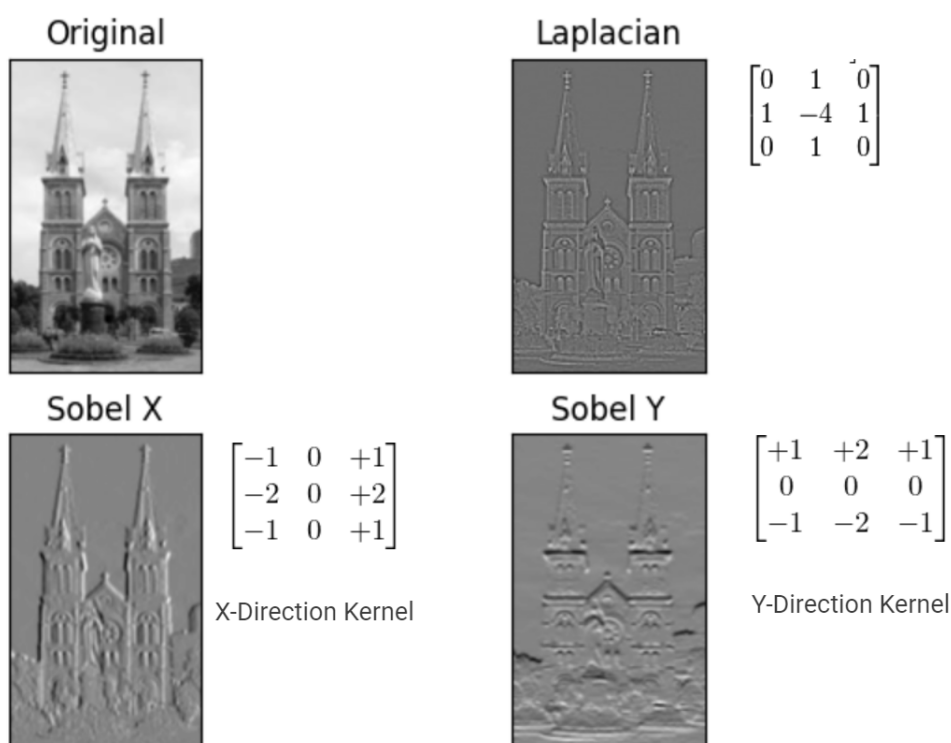
Sobel filters là edge-detectors, vertical Sobel filter sẽ tìm vertical edges, và horizontal Sobel filter tìm horizontal edges. Với những pixel sáng hơn, tức có giá trị cao hơn (gần 255) ở output image. Vậy Sobel filter cho ta cách filter cạnh theo chiều dọc và ngang của input

-1	0	1	1	2	1
-2	0	2	0	0	0
-1	0	1	-1	-2	-1

(a) Vertical Sobel filter

(b) Horizontal Sobel filter

Hình 11: Sobel filters



Hình 12: Sobel filter detect edges

image.

Tại sao edge-detected image mang thông tin hữu ích hơn ảnh gốc?

Khi CNN train với dataset MNIST, dựa vào edge-detection filter và so sánh các vertical edges gần trung tâm bức ảnh sẽ dễ nhận diện thuộc ký tự nào hơn. Nói chung, convolution giúp ta tìm đặc trưng của ảnh tại các vị trí cụ thể (ví dụ edges).

(c) **Padding**

Khi tích chập 4x4 input image với 3x3 filter được 2x2 output image, nhưng nếu muốn giữ nguyên kích thước output image như input image, ta phải thêm các giá trị vào ngoài đường bao của input image, gọi là padding. Giá trị thêm vào có thể là zero hoặc 255 hoặc một giá trị trung bình nào đó. Với 3x3 filter sẽ cần 1 pixel padding, gọi là **same padding** (sẽ có dạng khác gọi là valid padding):

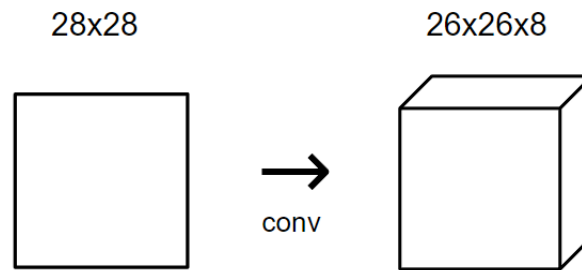
0	0	0	0	0	0
0	0	50	0	29	0
0	0	80	31	2	0
0	33	90	0	75	0
0	0	9	0	95	0
0	0	0	0	0	0

Hình 13: Same padding

(d) **Conv layer**

Từ các mục ở trên, ta đã biết image convolution hoạt động như thế nào và tại sao nó hữu dụng, tiếp theo sẽ tìm hiểu nó thực sự được dùng trong CNN ra sao. CNN bao gồm conv layer dùng các filter để xử lý các input images và cho ra output images. Tham số đầu tiên của conv layer là số lượng filter.

Trong MNIST CNN, chúng ta dùng một conv layer với 8 lớp filter như layer khởi tạo trong network, có nghĩa là 28x28 input image sẽ có 26x26x8 output volume (khối):



Hình 14: Conv layer

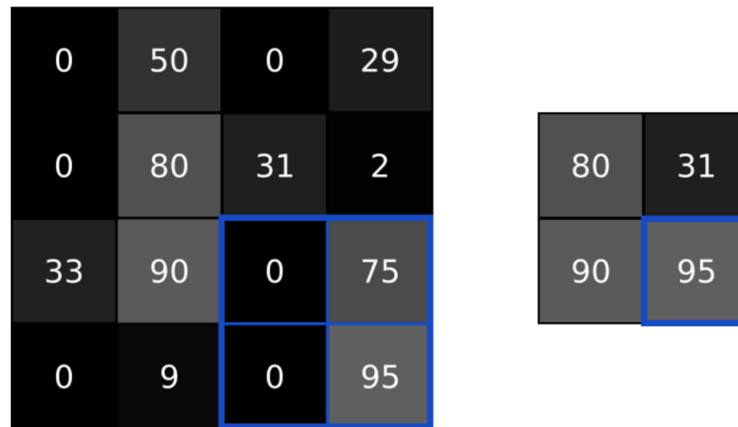
Ở đây output là 26x26x8 chứ không phải 28x28x8 vì ta dùng valid padding, nó sẽ giảm 2 pixels ở chiều ngang và chiều dọc. Ta có 3×3 (filter size) \times 8 (số filters) = 72 tham số.

(e) **Pooling**

Các pixels kế nhau trong 1 bức ảnh, có giá trị gần giống nhau, vì vậy conv layer sẽ cũng tạo ra các giá trị giống nhau cho các pixels kế cạnh trong output. Vì vậy thông tin chứa trong các output của conv layer là gần giống nhau, nên ta sẽ không có được thông tin mới.

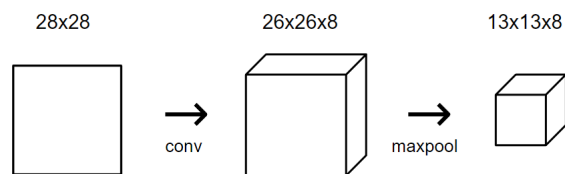
Pooling layer giải quyết được vấn đề này. Chúng giảm kích thước của input thông qua pooling values. Pooling thông thường được tính bằng min, max và average. Bên dưới là ví dụ về việc tính Max pooling với pooling size bằng 2:

Max pooling (pooling size = 2) trên 4x4 image sẽ giảm chiều output còn 2x2. Vậy mục tiêu của pooling là để thay đổi số lượng chiều dữ liệu, bằng cách chia chiều dọc và chiều ngang



Hình 15: Max Pooling with size 2x2

của input cho pool size.



Hình 16: Max pooling reduce dimentions of tensor

```

1 class MaxPool2:
2     # A Max Pooling layer using a pool size of 2.
3
4     def iterate_regions(self, image):
5         # Generates non-overlapping 2x2 image regions to pool over, image is a 2d
6         # numpy array
7
8         h, w, _ = image.shape
9         new_h = h // 2
10        new_w = w // 2
11
12        for i in range(new_h):
13            for j in range(new_w):
14                im_region = image[(i * 2):(i * 2 + 2), (j * 2):(j * 2 + 2)]
15                yield im_region, i, j
16
17    def forward(self, input):
18        # Performs a forward pass of the maxpool layer using the given input.
19        # Returns a 3d numpy array with dimensions (h / 2, w / 2, num_filters),
20        # input is a 3d numpy array with dimensions (h, w, num_filters)
21
22        self.last_input = input
23
24        h, w, num_filters = input.shape
25        output = np.zeros((h // 2, w // 2, num_filters))
26
27        for im_region, i, j in self.iterate_regions(input):
28            output[i, j] = np.amax(im_region, axis=(0, 1))

```

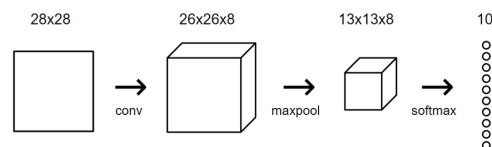
```

27
28     return output
29
30 def backprop(self, d_L_d_out):
31     # Performs a backward pass of the maxpool layer.
32     # Returns the loss gradient for this layer's inputs, d_L_d_out is the
    loss gradient for this layer's outputs.
33
34     d_L_d_input = np.zeros(self.last_input.shape)
35
36     for im_region, i, j in self.iterate_regions(self.last_input):
37         h, w, f = im_region.shape
38         amax = np.amax(im_region, axis=(0, 1))
39
40         for i2 in range(h):
41             for j2 in range(w):
42                 for f2 in range(f):
43                     # If this pixel was the max value, copy the gradient to it.
44                     if im_region[i2, j2, f2] == amax[f2]:
45                         d_L_d_input[i * 2 + i2, j * 2 + j2, f2] = d_L_d_out[i, j, f2]
46
47     return d_L_d_input
48

```

(f) Softmax

Để hoàn tất CNN, ta phải đưa ra khả năng dự đoán, dùng standard final layer cho multiclass classification: softmax layer, một fully-connected (dense) layer dùng softmax function như là activation.



Hình 17: Softmax

Ở một khía cạnh khác, có thể chuyển output thành xác suất, vậy số nào có xác suất dự đoán cao thì nó mang ký số đó, ta thực sự không cần dùng softmax để dự đoán các con số, mà chỉ đi tìm những số có output cao nhất từ network.

```

1 class Softmax:
2     # A standard fully-connected layer with softmax activation.
3
4     def __init__(self, input_len, nodes):
5         # We divide by input_len to reduce the variance of our initial values
6         self.weights = np.random.randn(input_len, nodes) / input_len
7         self.biases = np.zeros(nodes)
8
9     def forward(self, input):
10        # Performs a forward pass of the softmax layer using the given input.
11        # Returns a 1d numpy array containing the respective probability values,
        input can be any array with any dimensions.
12
13        self.last_input_shape = input.shape
14

```

```

15     input = input.flatten()
16     self.last_input = input
17
18     input_len, nodes = self.weights.shape
19
20     totals = np.dot(input, self.weights) + self.biases
21     self.last_totals = totals
22
23     exp = np.exp(totals)
24     return exp / np.sum(exp, axis=0)
25
26 def backprop(self, d_L_d_out, learn_rate):
27     # Performs a backward pass of the softmax layer.
28     # Returns the loss gradient for this layer's inputs. d_L_d_out is the
29     # loss gradient for this layer's outputs, learn_rate is a float
30     # We know only 1 element of d_L_d_out will be nonzero
31     for i, gradient in enumerate(d_L_d_out):
32         if gradient == 0:
33             continue
34
35         # e^totals
36         t_exp = np.exp(self.last_totals)
37
38         # Sum of all e^totals
39         S = np.sum(t_exp)
40
41         # Gradients of out[i] against totals
42         d_out_d_t = -t_exp[i] * t_exp / (S ** 2)
43         d_out_d_t[i] = t_exp[i] * (S - t_exp[i]) / (S ** 2)
44
45         # Gradients of totals against weights/biases/input
46         d_t_d_w = self.last_input
47         d_t_d_b = 1
48         d_t_d_inputs = self.weights
49
50         # Gradients of loss against totals
51         d_L_d_t = gradient * d_out_d_t
52
53         # Gradients of loss against weights/biases/input
54         d_L_d_w = d_t_d_w[np.newaxis].T @ d_L_d_t[np.newaxis]
55         d_L_d_b = d_L_d_t * d_t_d_b
56         d_L_d_inputs = d_t_d_inputs @ d_L_d_t
57
58         # Update weights / biases
59         self.weights -= learn_rate * d_L_d_w
60         self.biases -= learn_rate * d_L_d_b
61
62     return d_L_d_inputs.reshape(self.last_input_shape)

```

Phần II: Triển khai CNN bằng Python và Keras

Ta đi thực hiện CNN với dataset MNIST, dùng 3 lớp: Conv 3x3, MaxPool và Softmax.

1. import libraries và load dataset

```

1 from keras.datasets import mnist

```

```

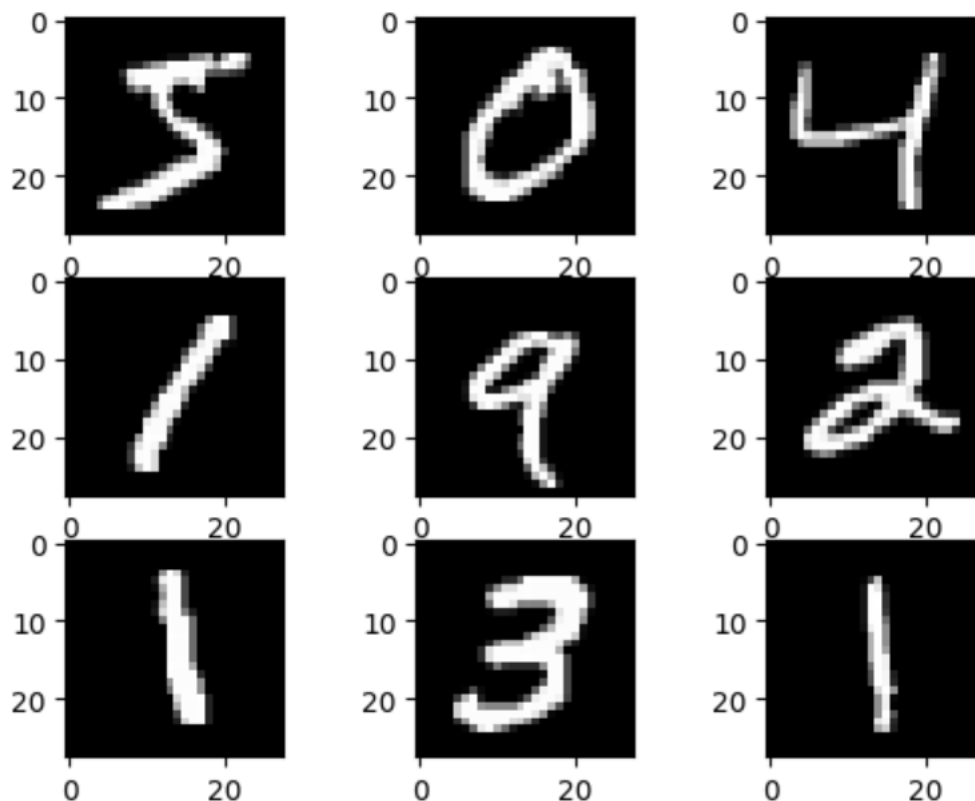
2 from matplotlib import pyplot
3 import numpy as np
4
5 #loading the dataset
6 (train_X, train_y), (test_X, test_y) = mnist.load_data()
7
8 # summarize loaded dataset
9 print('Train: X=%s, y=%s' % (train_X.shape, train_y.shape))
10 print('Test: X=%s, y=%s' % (test_X.shape, test_y.shape))
11 # plot first few images
12 for i in range(9):
13     # define subplot
14     pyplot.subplot(330 + 1 + i)
15     # plot raw pixel data
16     pyplot.imshow(train_X[i], cmap=pyplot.get_cmap('gray'))
17 # show the figure
18 pyplot.show()
19

```

Output:

Train: X=(60000, 28, 28), y=(60000,)

Test: X=(10000, 28, 28), y=(10000,)



Hình 18: cnn-implement-dataset

2. CNN with keras

```
1 from keras.models import Sequential
2 from keras.layers import Conv2D, MaxPooling2D, Dense, Flatten
3 from keras.utils import to_categorical
4 from keras.optimizers import SGD
5
6 learning_rate = 0.001
7 num_epochs = 30
8 batch_size = 1
9
10 train_images = train_X
11 train_labels = train_y
12 test_images = test_X
13 test_labels = test_y
14
15 train_images = (train_images / 255) - 0.5
16 test_images = (test_images / 255) - 0.5
17
18 train_images = np.expand_dims(train_images, axis=3)
19 test_images = np.expand_dims(test_images, axis=3)
20
21 model = Sequential([
22     Conv2D(8, 3, input_shape=(28, 28, 1), use_bias=False),
23     MaxPooling2D(pool_size=2),
24     Flatten(),
25     Dense(10, activation='softmax'),
26 ])
27
28 model.compile(SGD(learning_rate=learning_rate), loss='categorical_crossentropy',
29               metrics=['accuracy'])
30
31 history = model.fit(
32     train_images,
33     to_categorical(train_labels),
34     batch_size=batch_size,
35     epochs=num_epochs,
36     validation_data=(test_images, to_categorical(test_labels)),
37 )
```

Output:

```
1 #plot the loss and validation loss of the dataset
2 pyplot.plot(history.history['loss'], label='loss')
3 pyplot.plot(history.history['val_loss'], label='val_loss')
4 pyplot.legend()
5
```

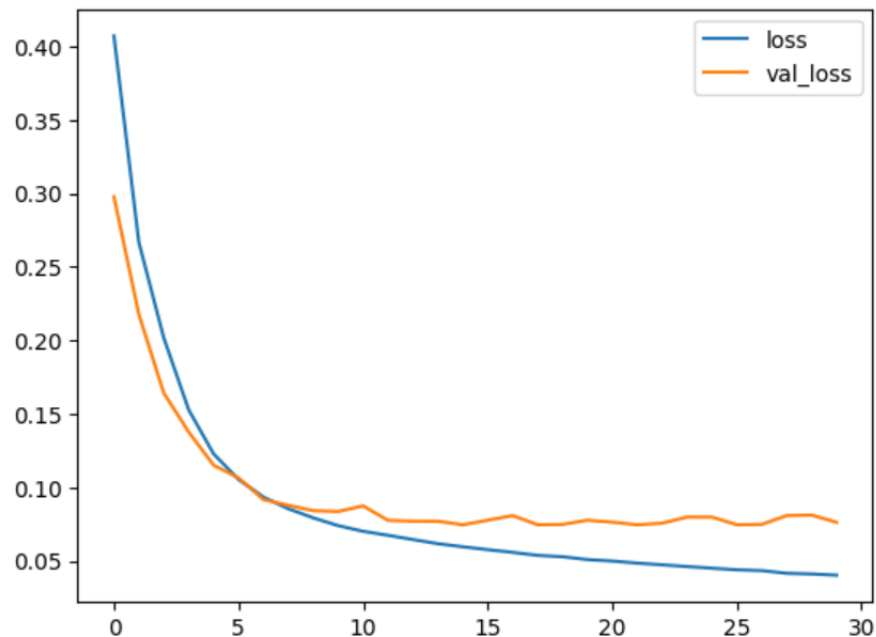
Output:

```

Epoch 7/30
60000/60000 [=====] - 98s 2ms/step - loss: 0.0933 - accuracy: 0.9738 - val_loss: 0.0917 - v
Epoch 8/30
60000/60000 [=====] - 102s 2ms/step - loss: 0.0854 - accuracy: 0.9758 - val_loss: 0.0877 - v
Epoch 9/30
60000/60000 [=====] - 88s 1ms/step - loss: 0.0792 - accuracy: 0.9768 - val_loss: 0.0841 - v
Epoch 10/30
60000/60000 [=====] - 82s 1ms/step - loss: 0.0739 - accuracy: 0.9786 - val_loss: 0.0836 - v
Epoch 11/30
60000/60000 [=====] - 66s 1ms/step - loss: 0.0701 - accuracy: 0.9795 - val_loss: 0.0874 - v
Epoch 12/30
60000/60000 [=====] - 85s 1ms/step - loss: 0.0673 - accuracy: 0.9808 - val_loss: 0.0775 - v
Epoch 13/30
...
Epoch 29/30
60000/60000 [=====] - 113s 2ms/step - loss: 0.0410 - accuracy: 0.9878 - val_loss: 0.0811 - v
Epoch 30/30
60000/60000 [=====] - 113s 2ms/step - loss: 0.0402 - accuracy: 0.9881 - val_loss: 0.0762 - v

```

Hình 19: Kết quả train model dùng keras



Hình 20: So sánh loss giữa train và validation