



# A comprehensive survey on convolutional neural network in medical image analysis

Xujing Yao<sup>1</sup> · Xinyue Wang<sup>1</sup> · Shui-Hua Wang<sup>1,2,3</sup> · Yu-Dong Zhang<sup>1,2</sup> 

Received: 21 May 2020 / Revised: 30 July 2020 / Accepted: 13 August 2020 /  
Published online: 24 August 2020

© Springer Science+Business Media, LLC, part of Springer Nature 2020

## Abstract

CNN is inspired from Primary Visual (V1) neurons. It is a typical deep learning technique and can help teach machine how to see and identify objects. In the most recent decade, deep learning develops rapidly and has been well used in various fields of expertise such as computer vision and natural language processing. As the representative algorithm of deep learning, Convolution Neural Network (CNN) has been regarded as a breakthrough of historic significance in image processing and visual recognition tasks since the astonishing results achieved on ImageNet Large Scale Visual Recognition Competition (ILSVRC) Unlike methods based on handcrafted features, CNN models can build high-level features from low-level ones in a data-driven fashion and have displayed great potential in medical image analysis among the aspects of segmentation of histological images identification, lesion detection, tissue classification, etc. This paper provides a review on CNN from the perspectives of its basic mechanism introduction, structure, typical architecture and main application in medical image analysis through analyzing over 100 references from Google Scholar, PubMed, Web of Science and various sources published from 1958 to 2020.

**Keywords** Deep learning · Feedforward Neural Network · Convolutional neural network · Breast Cancer · Lung Nodule · Brain Tumor · Medical image analysis

---

✉ Shui-Hua Wang  
shuihuawang@iecc.org

✉ Yu-Dong Zhang  
yudongzhang@iecc.org

Xujing Yao  
xy147@le.ac.uk

Xinyue Wang  
fireswxy@163.com

Extended author information available on the last page of the article

## 1 Background

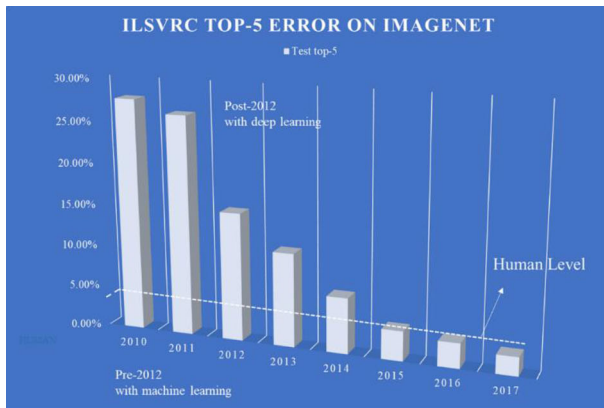
Deep Learning (DL) is a promising research direction in the Machine Learning (ML) area. By means of utilizing DL technique, computational models are capable to learn the features and representation of data with a process of extraction. This is an innovation that brings dramatical improvements in image recognition, object detection, voice recognition, genomic science and many other domains in real life [51, 84].

The study of machine learning is inspired in part by the exploration of human brain. Around 300 B.C., the concept of understanding the brain has attracted a lot of attention of researchers, among which the celebrated philosopher Aristotle is considered as the forerunner, due to the reason that he opens the history of human trying to learn the brain [72]. In 1943, the team of McCulloch and Pitts [61] brought the earlier neural network model: MCP Neural Model into the attention of public. About one decade later, Rosenblatt [75] introduced the perceptron based on this MCP model to complete the classification task. However, just one decade after the introduction, a significant flaw was identified by Minsky and Papert [64] that the perceptron cannot handle in the nonlinear classification problems, which hindered the subsequent progress of neural network for approximately twenty years [107]. In 1986, a book published by Rumelhart, Hinton [78] bring the turnaround for the development of neural network. The book detailedly introduced the back propagation algorithm as well as the whole derived process of utilizing back propagation algorithm and the activation function Sigmoid to deal with the nonlinear problems. This not only led to a great deal of uproar in academia but also attracted a sharp increase number of people into the research of deep learning. In 1997, the CNN model is presented [49, 63, 86]. In the same year, Hochreiter and Schmidhuber [29] discovered Long Short-Term Memory (LSTM) Networks model. After the introduction of ReLU, neural network has so far been known as the most outstanding tool to deal with learning problems [1, 8].

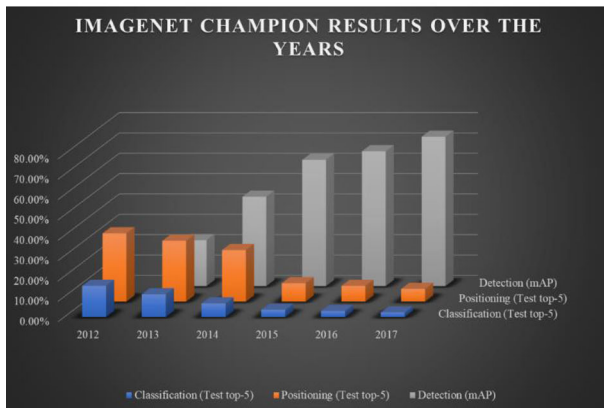
Because of the efforts of so many great scholars, as shown in Fig. 1, deep learning is becoming more and more accurate and has evolved to perform many complex tasks and achieve loads of unimaginable breakthroughs [22, 30, 48, 95].

This survey will introduce the application of CNN technology in analysing medical images. The reason why we chose to introduce CNN is that it can process image information better than deep neural network (DNN). DNN is one of the simplest neural networks. In DNN, each neuron belongs to a different layer and is connected with all neurons in the previous layer. Signal is transmitted one-way from the input layer to the output layer. While CNN is a feedforward neural network through convolution computation, which is proposed by the biological receptive field mechanism and has translational invariance. It utilizes the convolution kernel to maximize the application of local information and retain the planar structure information. DNN is input in vector form, without considering the structure information of the plane. In the field of image processing, the plane information is very significant, so CNN performs better than DNN in medical image analysis [55].

A simple flowchart of this survey is organized as Fig. 2. The purpose and contribution of this review is to: (i) introduce the theoretical basis and working principle of CNN in detail; (ii) elaborate description of the popular forms of various CNN basic structures; (iii) detailedly explain classic models and trends from the rise of CNN to the present. (iv) demonstrate the outstanding contributions of CNN models in the medical field from recent years. It can be found from our review that CNN model is increasingly being used in medical fields where it is difficult to find the basis for diagnosis, or hard to achieve high accuracy or fast diagnostic analysis by just utilizing human power. Perhaps in the next few years, humans will be able to



(a) ILSVRC Top-5 Error on ImageNet



(b) ImageNet Champion results over the years

Fig. 1 Breakthroughs in ImageNet competitions over recent years

use CNN to solve medical problems that are currently difficult to understand and analyze from a medical perspective.

Section 2, 3 and 4 will explain CNN from three aspects: the basic operating mechanism, the basic structure and typical architectures of CNN. The operating mechanism section introduces the operation unit, operation thought, gradient vanishing, feedforward propagation and feedback propagation. The basic structure section describes principles and details of convolutional layer, pooling layer and fully connected layer. The classic architecture section reviews the theory and advantages of several typical architectures such as AlexNet, VGG-Net, Inception, ResNet, R-CNN, SPP-Net and so on.

## 2 Basic operating mechanism

### 2.1 Operation unit: Neuron

In 1943, the MCP Model introduced by McCulloch and Pitts [61] is very similar with biological neural network. From a biological kind of view, every neuron is linked with other

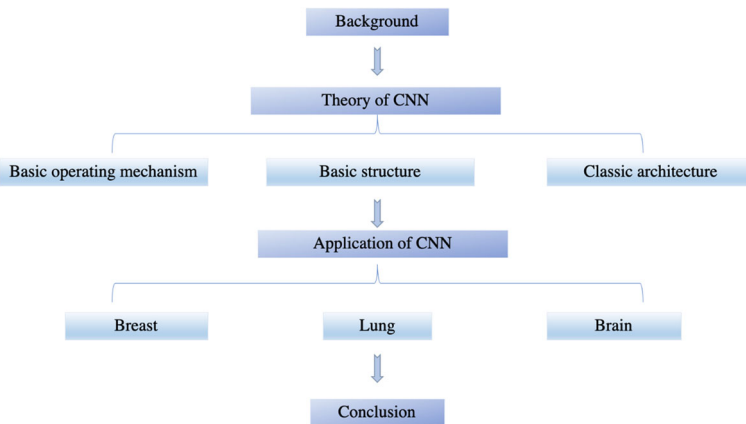


Fig. 2 Simple flowchart of this survey

neurons. When the neuron is activated, which means the potential of neuron is over a threshold, it will send chemical substances to the connected neurons to transform the potential of other neurons. MCP model has the same mechanism. The circle in Fig. 3 represents neuron, which receives the input signals  $(x_1 \dots x_m)$  from the other neurons. These input signals will get different weights  $(w_1 \dots w_m)$  and bias  $(b)$  and then go to the neuron. The output  $\hat{y}$  will be got through the activation function [1, 59].

$$\hat{y} = \sigma(z) = \sigma\left(\sum_{j=1}^m w_j x_j + b\right) \tag{1}$$

Where we assume the input of the activation function as  $z$ .

In a multi-layer neural network, each neuron node will accept the output of previous layer as its input value and then pass the value to the following layer. This function relationship between neuron nodes is defined as the activation function. If you just utilize the linear

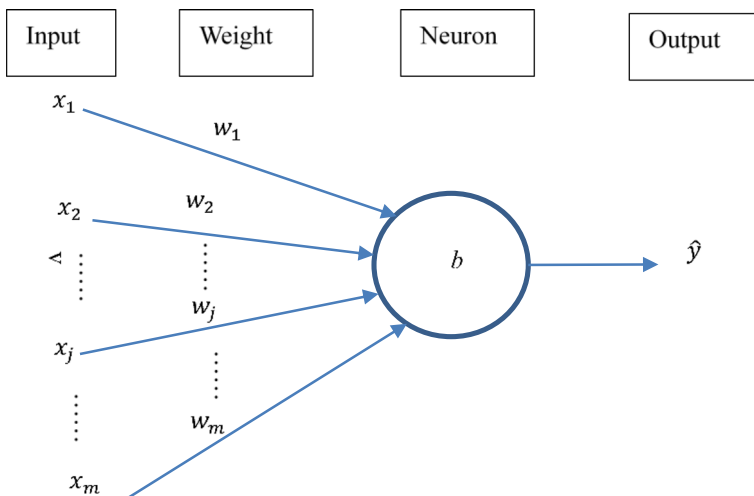


Fig. 3 The McCulloch and Pitts model of a neuron

activation function or do not use the activation function, all things are about calculating the linear activation function, which is very limited, so in order to get more possibilities, the nonlinear activation function is born [25, 67, 82].

## 2.2 Activation function and loss function

### 2.2.1 Sigmoid function

Sigmoid Function aims to alter the continuous real value of the input  $z$  to the output  $\sigma(z)$  of a range between 0 and 1.

$$\sigma_{\text{Sigmoid}}(z) = \frac{1}{1 + e^{-z}} \tag{2}$$

It's clear from Eq. (2) that if the input is an extremely large negative number, the output is going to be 0. If the input value is an extremely large positive number, the output will be 1. But this function also has obvious disadvantages that it consumes loads of calculation time, and will be easy to produce gradient disappearance and gradient explosion [16]. In addition, as displayed in Fig. 4a, its data is not centralized and symmetric with (0,0.5) as the center, which will bring a negative effect on the gradient calculation of the hidden layer.

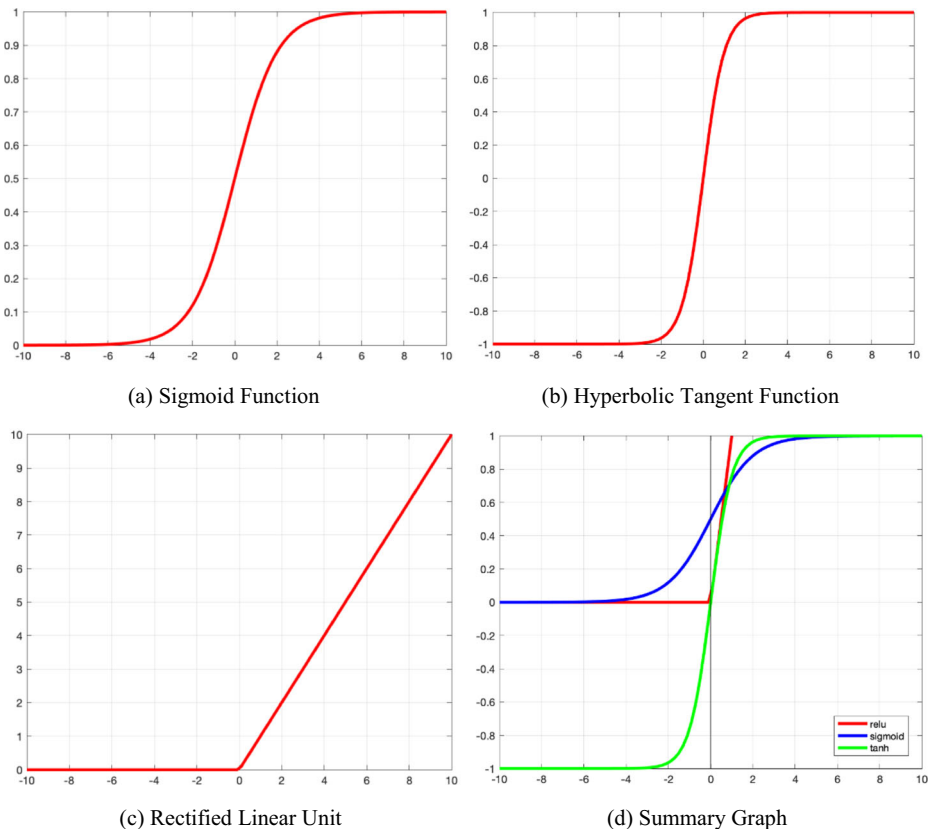


Fig. 4 Activation functions

## 2.2.2 Hyperbolic tangent function (Tanh)

Tanh function performs superior to the sigmoid function in hidden layer. It can be viewed as the shifted version of sigmoid as shown in Fig. 4b. The equation of Tanh is shown as:

$$\sigma_{\text{Tanh}}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (3)$$

Although it achieves a similar effect of data centralization, in practice, sigmoid and tanh functions tend to reach saturation at the end value, especially in the deep network training, which leads to a slow training speed.

## 2.2.3 Rectified linear unit function (ReLU)

It can be seen from Fig. 4c, ReLU solves the gradient disappearance problem when there is a positive input [6]. The equation is shown below

$$\sigma_{\text{ReLU}}(z) = \max(0, z) \quad (4)$$

Moreover, ReLU costs much less computation time when compared with sigmoid and tanh function. But the gradient vanishing problem remains and the data distribution of ReLU is not centralized when there is a negative input [69].

## 2.2.4 Loss function

After the activation function, the calculated value  $\hat{y}_j$  is got. However, the calculated value cannot equal to the actual value  $y_j$  due to some unexplained factors. The loss function aims to minimize the distance between each actual value and the predicted value in the training set, which refers to, the sum of the error  $J(w)$  between  $\hat{y}_j$  and  $y_j$  through utilizing the mean square error principle shown below:

$$J(w) = \frac{1}{2} \sum_{j=1}^m (\hat{y}_j - y_j)^2 \quad (5)$$

Where  $w$  represents weight, which is the coefficient of linear relation.

However, the mean square error has many disadvantages when it is used in the back propagation algorithm which will be introduced in Section 2.4.

## 2.3 Operation thought: Gradient descent

Most deep learning algorithms involve many kinds of optimizations. Gradient descent method is one of the optimizations which used to minimize the loss function. In this method, the minimum loss function and the value of model parameters are obtained by step-by-step iteration [25].

In 1847, Cauchy, Augustin discovered gradient descent. Assume there is a function:

$$y = f(x) \quad (6)$$

Where both  $x$  and  $y$  are real numbers.

The derivative of the function is  $[dy/dx]$  or  $f'(x)$ , which presents the slope of  $f(x)$  at point  $x$ . Besides, it also shows if the input has a little change like  $\epsilon$ , the output can also be represented by the slope and the original output, the formula is shown below [52]:

$$f(x) \leftarrow f(x + \epsilon) = f(x) + \epsilon f'(x) \tag{7}$$

Where  $\leftarrow$  represents the assignment operation.

With this equation, it is easy to control the output decrease. When the slope is negative the input should become bigger, and when the slope is positive the input should become smaller in order to get a smaller output value as the Fig. 5 shows.

Since the loss function should be smallest, the equation will run until it meets the three possibilities as shown in Fig. 6. First, the output is the minimum value within limits, which is called local minimum. Second, the output is the lowest point in the whole function, which is called global minimum. And the third possibility is that the output is a saddle point. A saddle point refers to a critical point which is neither maximum nor minimum.

In order to avoid the output which comes from the saddle point, it is important to choose the suitable learning rate  $\epsilon$ . As displayed in Fig. 7, if the learning rate is chosen suitable, the minimum value can be found easily as the orange line. But if  $\epsilon$  is too small as the green line, the time to run the equation will increase and the output may be generated from a saddle point or local minimum. If  $\epsilon$  is large as the blue line, it is impossible to find the minimum. If  $\epsilon$  is too large as the red line, the equation will loss control and the output is far from the right point. Therefore, it is important to control learning rate  $\epsilon$ . It is convenient and visualized to draw a graph about loss and parameters change. And from Fig. 7, it is obvious that the orange line gives us the best learning rate.

Except choosing the suitable learning rate, different kinds of gradient descent can improve the loss function in different levels. We assume the original equation is

$$\hat{y} = b + w_1x_1 + w_2x_2 + \dots w_mx_m \tag{8}$$

The following introductions and equations come from the loss function based on mean squared error:

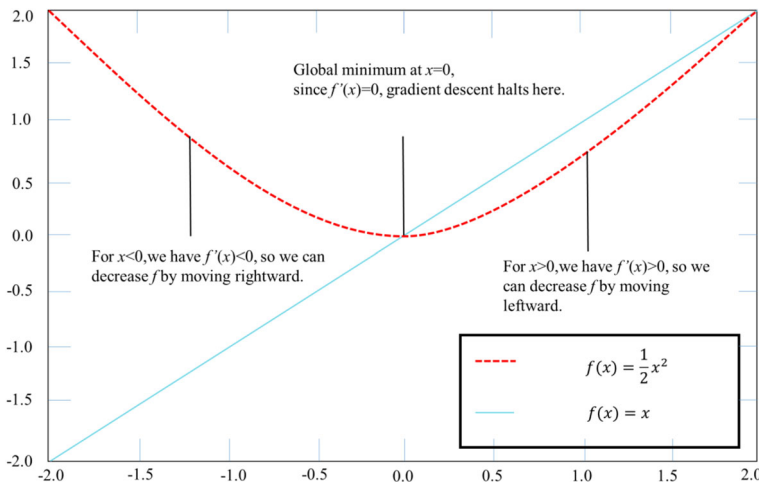


Fig. 5 The relationship between a function and the derivatives of a function

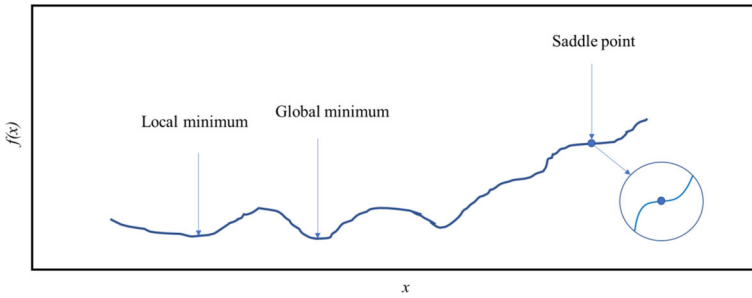


Fig. 6 Three possibilities the formula will stop when it meets

$$J(w) = \frac{1}{2} \sum_{l=1}^n \sum_{j=1}^m (\hat{y}(x_j^l) - y^l)^2 \tag{9}$$

Where:

$x_j^l$  the  $j$ -th neuron in the  $l$ -th layer.

$n$  the total number of layers.

Three basic gradient descent variants are introduced next.

### 2.3.1 Batch gradient descent (BGD)

As shown in Eq. (10), BGD requires the calculation of partial derivatives for the whole data set according to Eqs. (7) and (9). To get the optimum  $w$  for achieving the maximum value of  $J(w)$ , we will repeat the following equation until convergence:

$$w_i \leftarrow w_i + \epsilon \sum_{l=1}^n \sum_{j=1}^m (y^l - \hat{y}(x_j^l)) x_{i,j}^l \tag{10}$$

Where  $\leftarrow$  represents the assignment operation.

$\epsilon$  represents the learning rate.

$i$  represents the number of iterations.

However, it also has obvious disadvantages. In one update, it calculates the gradient for the whole data set, so the calculation time is very long. In this way, it will be very difficult to deal with a large scale data set, and update the model in real time with new data [77].

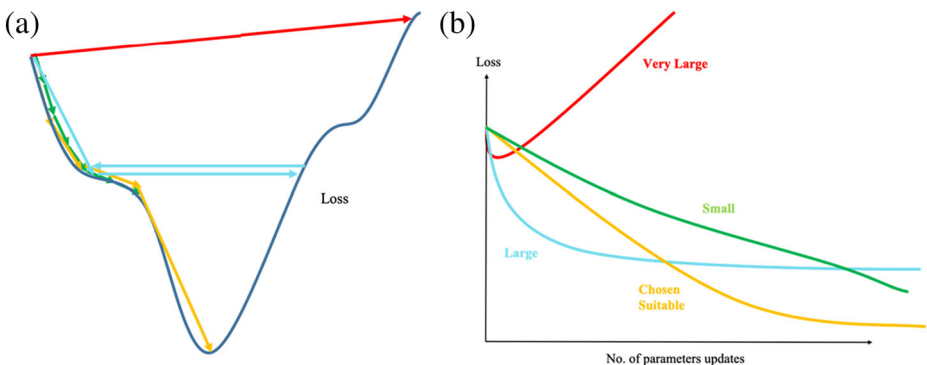


Fig. 7 The relationship between loss and parameters change



### 2.3.2 Stochastic gradient descent (SGD)

For a large data set, there may be similar samples, so the BGD method of calculating the gradient with all the data at one time will be redundant. While SGD will update each sample with a gradient every time, which is not redundant and relatively fast, and new samples can be added in. Its equation is shown in Eq. (11). However, for accuracy, SGD only uses one sample to determine the direction of the gradient, resulting in a solution that is probably not optimal. However, due to the reason that SGD is updated frequently, it may lead to serious oscillation of cost function [10, 11].

$$w_i \leftarrow w_i + \epsilon \left( y^j - \hat{y}(x_j^l) \right) x_{i,j}^l \quad \forall j = 1, \dots, m \wedge \forall l = 1, \dots, n \tag{11}$$

### 2.3.3 Mini-batch gradient descent (MBGD)

The difference between MBGD and SGD is that MBGD used a small batch of sample data consists of  $n$  samples to calculate the loss function during each update. The value range of the hyperparameter  $n$  is generally between 50 and 256. In this method, as shown in Eq. (12), when updating the parameters, the variance will be reduced, the convergence will be more stable, and the calculation speed will be faster. However, if  $\epsilon$  is extremely small, the convergence rate will be slow. Besides, if the loss is extremely large, the loss function will oscillate near the minimum value. In addition, the choice of batch size is also difficult [28, 47].

$$w_i \leftarrow w_i - \epsilon \sum_{l=t}^{t+n-1} \sum_{j=1}^m \left( \hat{y}(x_j^l) - y^j \right) x_{i,j}^l \tag{12}$$

### 2.3.4 Gradient vanishing and gradient explosion

Gradient disappearance and gradient explosion are essentially the same, both of which are due to the continuous multiplication effect in gradient back propagation brought by the too deep network layers. In other words, if the numbers in each layer are smaller than 1, then the positive gradient that you get when you multiply them is very small, and that's the gradient disappearing. On the contrary, if the numbers in each layer are bigger than 1, then the forward gradient will be extremely large, that is, the gradient explosion. There are many ways to solve these two problems, such as switching to activation functions such as ReLU, LeakyReLU, ELU, utilizing batch normalization, ResNet residual structure, LSTM structure, pretraining plus finetuning, etc. [9, 68].

## 2.4 Propagation thought

### 2.4.1 Feedforward propagation

Like the introduced in Section 2.1, The core idea of MCP is a linear relationship plus an activation function in a small local model, which is:

$$y = \sigma(z) = \sigma_{\text{ReLU}}(b + \sum w * x). \tag{13}$$

Where  $*$  represents the convolution,  $b$  represents bias,  $\sigma$  represents activation function, usually refers to ReLU.

And since MCP is very simple which can only solve binary classification, the Multi-Layer perceptron (MLP) is introduced which is also called deep neural network (DNN). The major difference between DNN and MCP is the number of layers. With the core idea of MCP as the feedforward propagation of DNN, the general formula is:

Assume  $l-1$ -th layer has  $m$  neurons:

$$y_j^l = \sigma_{\text{ReLU}}(z_j^l) = \sigma_{\text{ReLU}}\left(\sum_{k=1}^m W_{jk}^l y_k^{l-1} + b_j^l\right) \quad (14)$$

where  $w_{jk}^l$  refers to weight between the  $j^{\text{th}}$  neurons in the  $l^{\text{th}}$  layer and the  $k^{\text{th}}$  neurons in the layer before.

It can be simplified that the output of  $l$ -th layer is:

$$y^l = \sigma_{\text{ReLU}}(z^l) = \sigma_{\text{ReLU}}\left(W^l y^{l-1} + b^l\right) \quad (15)$$

Where  $W^l$  refers to a  $n \times m$  matrix formed by linear coefficient  $w$  of layer  $l$ .

$y^{l-1}$  refers to a  $m \times 1$  vector formed by output  $y$  of layer  $l-1$ .

$b^l$  refers to a  $n \times 1$  vector formed by bias  $b$  of layer  $l$ .

$z^l$  refers to a  $n \times 1$  vector formed by linear output before activating of layer  $l$ .

However, DNN has a critical defect that if the input is large, a huge amount of weights and other parameters require being trained and is likely to lead to overfitting and local optimization. Therefore, the convolutional neural network (CNN) is discovered by LeCun, Boser [50] in 1989. It has two other kinds of layer, one is convolutional layer, the other is pooling layer. The detail information will introduce in Section 3.

## 2.4.2 Feedback propagation

According to the Section 2.1, loss function aims to minimize the difference between the calculated output ( $\hat{y}_j$ ) and the real output ( $y_j$ ). And the difference between them is caused by the random weights and bias. In order to get suitable weights and bias to make the calculated output correctly, the feedback propagation is used. The main concept of feedback propagation is to optimize the loss function by gradient descent. We assume  $W$  and  $X$  represent the weights matrix and input vector in one layer separately.

Since the loss function is

$$J(W) = \frac{1}{2} \sum_{l=1}^n (\hat{y}^l - y)^2 \quad (16)$$

Where:

$\hat{y}^l$  refers to calculated output in the  $l$ -th layer.

$y$  refers to real output in the  $l$ -th layer.

$n$  refers to total number of layers

$$\hat{y}^l = \sigma_{\text{ReLU}}(z^l) = \sigma_{\text{ReLU}}(b^l + \sum W^l * X^{l-1}) \quad (17)$$

The loss function's partial derivative respect to  $W^l$  is

$$\frac{\partial J(W)}{\partial W^l} = \left[ (\hat{y}^l - y) \odot \sigma'_{\text{ReLU}}(b^l + \sum W^l * X^{l-1}) \right] \left( \hat{y}^{l-1} \right)^T \quad (18)$$

and its partial derivative respect to  $b^l$  is

$$\frac{\partial J(W)}{\partial b^l} = \delta^l \left[ \left( \hat{y}^l - y \right) \odot \sigma'_{ReLU} \left( b^l + \sum W^l * X^{l-1} \right) \right] \tag{19}$$

Where  $\odot$  represents Hadamard Product, which is defined as the product of the corresponding elements of two matrices.

With the two Partial derivative formula, we utilize gradient descent to get the new weight  $W^l$  and bias of  $l$ -th layer  $b^l$  [13]. The formulas are as following:

$$W^l \leftarrow W^l - \varepsilon \frac{\partial J(W)}{\partial W^l} \tag{20}$$

$$b^l \leftarrow b^l - \varepsilon \frac{\partial J(W)}{\partial b^l} \tag{21}$$

Where  $\varepsilon$  is learning rate.

Assume the layer number of the last layer which is output layer is  $L$ . Then, the equation of gradient expression which is recorded as  $\delta$  of the  $L$  layer is shown below:

$$\delta^L = \frac{\partial J(W, b)}{\partial z^L} = \left( \hat{y}^L - y \right) \odot \sigma'_{ReLU} \left( b^L + \sum W^L * X^{L-1} \right) \tag{22}$$

With the above formulas, the gradient expression of layer can be calculated by the fully connected layer (FC layer) after it. Assuming the parameters and gradient expression of  $l$ -th fully connected layer are known, the gradient expression of  $(l - 1)$ -th layer can be got by this formula:

$$\delta^{l-1} = \frac{\partial J(W, b)}{\partial z^{l-1}} = \left( \frac{\partial z^l}{\partial z^{l-1}} \right)^T \frac{\partial J(W, b)}{\partial z^l} = \left( \frac{\partial z^l}{\partial z^{l-1}} \right)^T \delta^l \tag{23}$$

The relationship between  $z^l$  and  $z^{l-1}$  is

$$z^l = \hat{y}^{l-1} * W^l + b^l = \sigma_{ReLU} \left( z^{l-1} \right) * W^l + b^l \tag{24}$$

The above formula about the FC layer can be shown as following:

$$\delta^l = \left( \frac{\partial z^{l+1}}{\partial z^l} \right)^T \frac{\partial J(W, b)}{\partial z^{l+1}} = \left( W^{l+1} \right)^T \delta^{l+1} \odot \sigma'_{ReLU} \left( z^l \right) \tag{25}$$

When the  $l$ -th layer is convolutional layer, the gradient expression of  $l - 1$  layer which is hidden layer can be calculated by the formula shown below:

$$\delta^{l-1} = \left( \frac{\partial z^l}{\partial z^{l-1}} \right)^T \delta^l = \delta^l * \text{rot180} \left( W^l \right) \odot \sigma'_{ReLU} \left( z^{l-1} \right) \tag{26}$$

The convolution kernel has been rotated 180 degrees. So  $\text{rot180}()$  means flip up and down once, then flip left and right once. Since CNN also includes the pooling layer which compresses data, the data should be restored.

When the  $l$ -th layer is pooling layer, the calculation of gradient expression of  $(l - 1)$ -th layer is according to equation as:

$$\delta^{l-1} = \text{upsample}(\delta^l) \odot \sigma'_{ReLU}(z^{l-1}) \quad (27)$$

This up-sample function helps accomplishing the logic of pooling error matrix amplification and error redistribution [85].

In conclusion, with the basic operating mechanism about CNN, it is relatively easy to understand the structure of the CNN introduced in the following section.

### 3 Basic structure

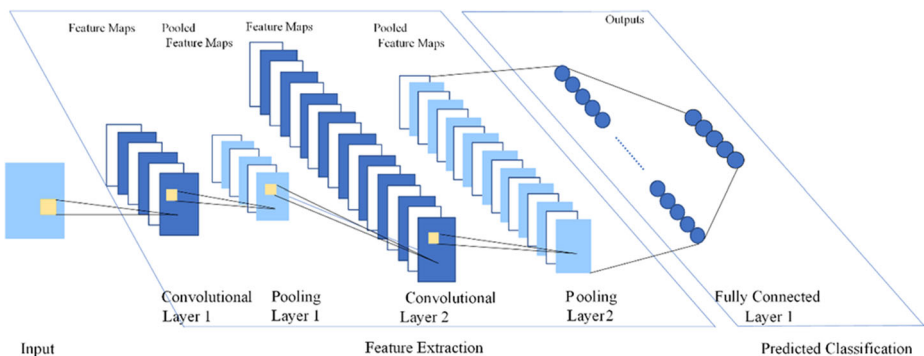
CNN can optimize the performance and generalization capability of the network by using multilayer structure [71, 76] illustrated in Fig. 8. The most commonly used layers are introduced in the next.

#### 3.1 Convolutional layer

The main ability of the convolutional layer is to extract critical features for getting a good prediction from the input images, which means reducing images into forms that are easier to process [17].

##### 3.1.1 Standard convolution

The element utilized in conducting the convolution operation in a convolutional layer is named the kernels whose parameters needs to be studied. The size of the kernel should be smaller than those of the input image. In the convolution process, each kernel will be convolved with the input image to compute an activation map. That is to say, the kernel will slide across the width and height of the input, calculating the dot product between the input and the kernel at each spatial position. Next, by superimposing the activation maps of each kernel along the depth dimension, the output image of the convolution layer can be achieved. It is because each kernel is smaller in width and height when compared to the input, each neuron in the activation map



**Fig. 8** A sample structure of CNN

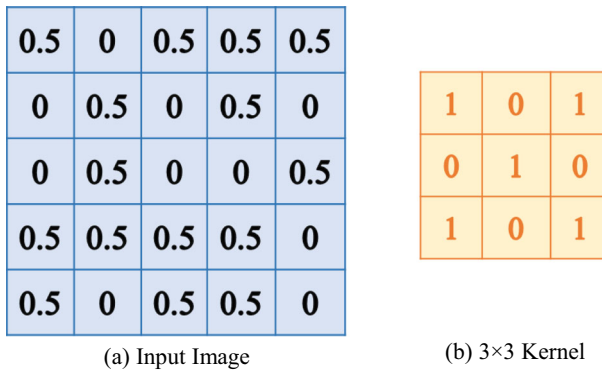


Fig. 9 Basic component of Standard Convolutional Layer

is associated only to a small local area of the input image, which means size of each neuron’s receptive field is small, equals to size of kernel [43].

Figure 9 is considering a  $5 \times 5$  input image in which matrix is composed of 0.5 and 0 and a  $3 \times 3$  kernel composed of 1 and 0. The size of kernel determines the receptive field of the convolution. The most preference of selecting kernel size for a 2D input image is 3, which means  $3 \times 3$  pixels.

Here the kernel will convolve over each cell position of the input image and eventually we will get a convoluted image. For every single point in it, you can take the point and the  $3 \times 3$  points around it out and do the convolution. For instance, for the lower right (row 4, column 4) point, the local convolution of the  $5 \times 5$  image and the  $3 \times 3$  matrix could be computed as shown in Fig. 10.

### 3.1.2 Dilated convolution

One disadvantage of standard convolution is its design for up-sampling and pooling. During these processes, there will be internal data structure loss and spatial hierarchical information loss. To deal with this loss problem, one significant hyper-parameter called dilation rate is introduced to convolutional layers. The dilation rate of kernel refers to the spacing between values inside (for a standard convolution, the dilation rate equals to 1) [56]. Through the utilizing of dilated convolution, the receptive field of a  $3 \times 3$  kernel with a dilation rate of 2 and no padding would as large as a  $5 \times 5$  kernel as illustrated in Fig. 11.

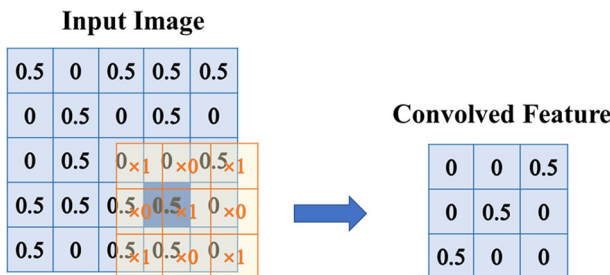
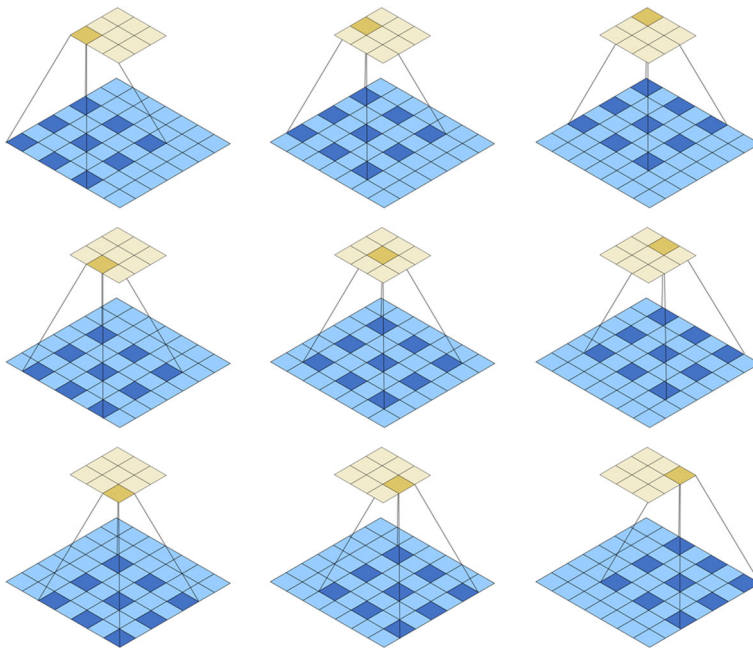


Fig. 10  $3 \times 3$  Convolution process



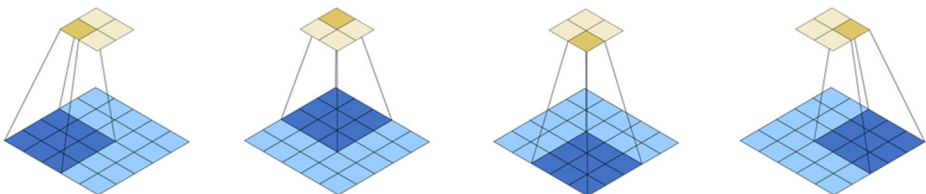
**Fig. 11** The process of convolution utilizing a  $3 \times 3$  kernel with a dilation rate of 2

In this way, a wider reception field would be provided at the same cost of computation which makes dilated convolution a promising solution among the area of real-time segmentation. It can be of great use especially when there is a need of wider field of view but larger kernel size or multiple convolutions cannot be afforded [100].

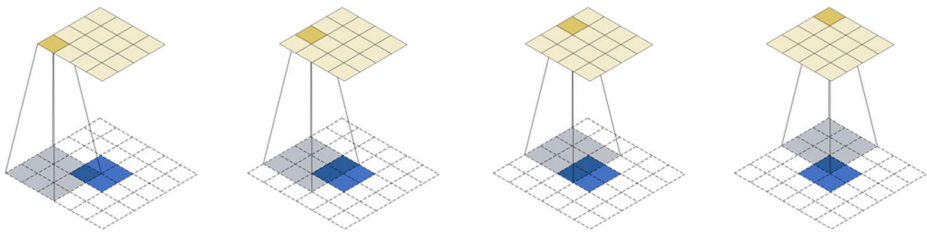
### 3.1.3 Transposed convolution

For a conventional convolution layer, forward propagation is ‘standard convolution’, which maps the input feature map into the output feature map as shown in Fig. 12, while reverse propagation is ‘transposed convolution’ [105, 106], which calculates the gradient of the input feature map according to the gradient of the output feature map and keeps the size of the gradient map same as that of the feature map illustrated in Fig. 13.

The process of transposed convolution is to some degree resemble to the deconvolution in that it produces the same spatial resolution as a deconvolutional layer. But the practical mathematical process carried out on the values is not the same. The transposed convolution



**Fig. 12** Part of the process of convolution utilizing a  $3 \times 3$  kernel over a  $5 \times 5$  input using unit-stride



**Fig. 13** Part of the transpose process of convolution utilizing a  $3 \times 3$  kernel over a  $4 \times 4$  input utilizing unit-stride (same to convolution utilizing a  $3 \times 3$  kernel over a  $2 \times 2$  input with a zero padding of 2 with unit-stride)

guarantees that the output will be as large as before, while it is still carrying out a standard convolution process. To accomplish the goal, some fancy padding on the input is performed.

The transposed convolution step will not reverse the operation just mentioned, especially not concerning the numeric values. It only reconstructs the spatial resolution and carries out a convolution. This is very helpful for Encoder-Decoder architectures. Through utilizing this, upscaling of an image with a convolution can be achieved at one time, rather than two separate operations.

To conclude, if convolution is regarded as the mapping of the map of input size to the map of output size, transposed convolution is the mapping of the map of output size to the map of input size [20], which aims to restore the size. Some typical cases of transposed convolution are displayed in Fig. 14.

### 3.1.4 Separable convolution

In neural networks, separable convolution is commonly used in order to reduce the computation. The convolution process can be divided into two steps: one depth wise convolution and one point wise convolution [15]. In the depth wise convolution process, we use three convolution kernels of  $5 \times 5 \times 1$  to produce an  $8 \times 8 \times 3$  image as illustrated in Fig. 15.

In the pointwise convolution process, first, the convolution of  $1 \times 1 \times 3$  is used to check the calculation of each pixel, and the image of  $8 \times 8 \times 1$  is obtained. Then, 256 convolution kernels of  $1 \times 1 \times 3$  are used to produce the image of  $8 \times 8 \times 256$  as displayed in Fig. 16.

Considering the change in computation, originally, 256 convolution kernels of  $5 \times 5 \times 3$  were moved  $8 \times 8$  times, that is,

$$256 \times 5 \times 5 \times 3 \times 8 \times 8 = 1,228,800 \tag{28}$$

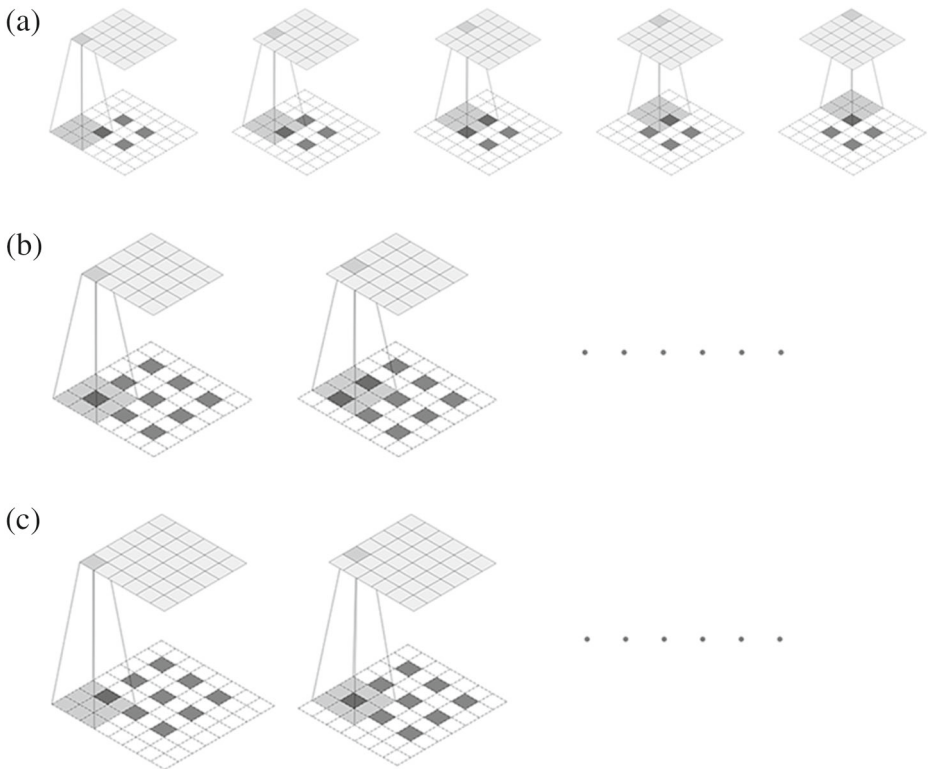
times of multiplication.

During the utilization process of depth wise convolution, 3 convolution kernels of  $5 \times 5 \times 1$  move  $8 \times 8$  times, requiring,

$$3 \times 5 \times 5 \times 7 \times 8 \times 8 = 4800 \tag{29}$$

times of multiplication.

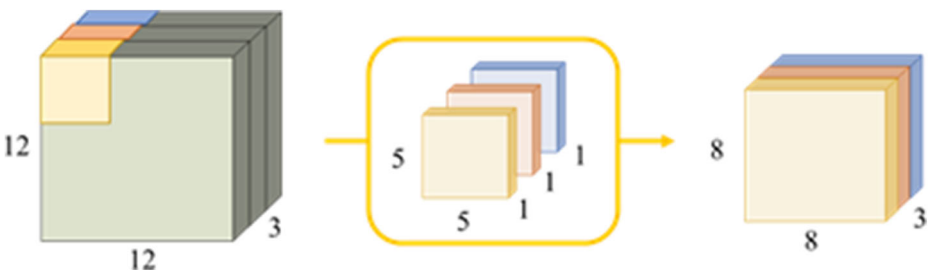
During the utilization process of pointwise convolution, 256 convolution kernels of  $1 \times 1 \times 3$  move  $8 \times 8$  times, requiring,



**Fig. 14** Typical cases of transposed convolution (a) Part of the transpose of convolution utilizing a  $3 \times 3$  kernel over a  $5 \times 5$  input utilizing  $2 \times 2$  strides (same to convolution utilizing a  $3 \times 3$  kernel over a  $2 \times 2$  input (gapped) with a zero padding of 2 utilizing unit-stride) (b) Part of the transpose of convolution utilizing a  $3 \times 3$  kernel over a  $5 \times 5$  input with a zero padding of 1 utilizing  $2 \times 2$  strides (same to convolution utilizing a  $3 \times 3$  kernel over a  $3 \times 3$  input (gapped) with a zero padding of 1 utilizing unit-stride) (c) Part of the transpose of convolution utilizing a  $3 \times 3$  kernel over a  $6 \times 6$  input with a zero padding of 1 utilizing  $2 \times 2$  strides (same to convolution utilizing a  $3 \times 3$  kernel over a  $3 \times 3$  input (gapped) with a zero padding of 1 utilizing unit-stride)

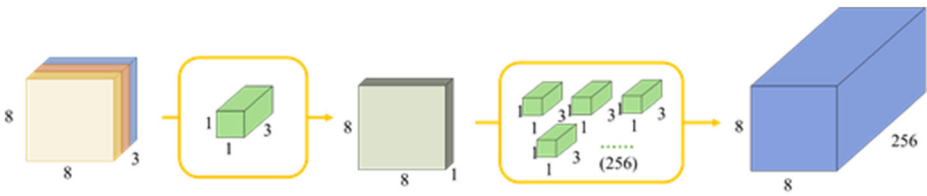
$$256 \times 1 \times 1 \times 3 \times 8 \times 8 = 49,152 \tag{30}$$

times of multiplication calculation, which adds up to 53,952 times of calculation Thus, the complete separable convolution process only needs 53,952 times of calculation which is much smaller than the computation of original process [14]. Separable convolution helps a lot in the application of mobile devices due to it uses the parameters efficiently [88].



**Fig. 15** Depth wise convolution: transform a  $12 \times 12 \times 3$  image to an  $8 \times 8 \times 3$  image by using 3 kernels





**Fig. 16** Pointwise convolution: transform an image from 3 channels to 1 channel, then output an image with 256 channels by using 256 kernels

### 3.2 Pooling layer

The function of pooling layer is reducing the resolution of feature maps. In other words, it aims to reduce the size of matrix generated by a convolutional layer. Typical categories of pooling will be introduced in the next [91].

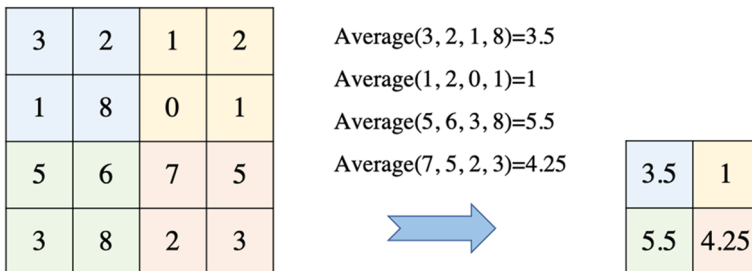
#### 3.2.1 Average pooling

In a simple term, average pooling means taking the average of a small range. Figure 17 indicates an instance of average pooling utilizing a kernel of  $2 \times 2$  with a stride of 2. After pooling, the average value in the  $2 \times 2$  kernel will be calculated and reserved in a rectified feature map. In this way, the  $4 \times 4$  input is compressed to  $2 \times 2$ .

Error of feature extraction mainly caused by two sources. On the one hand, the increasing variance of the estimated value led from the limitation of neighborhood size. On the other hand, the parameter error from the convolutional layer will result in the deviation of the estimated mean value. Generally, average pooling to some degree solve the first problem by retaining more background information of the input image. Average pooling emphasizes the sub-layer sampling of the overall feature information, reducing the parameter dimension while ensuring the integrity of the information transferred. In a large model such as Dense Net, average pooling is frequently utilized to shrink the dimension and at the same time facilitate the transfer of information to the following module for extracting features [108].

#### 3.2.2 Max pooling

In a simple term, max pooling means taking the maximum of a small range. Figure 18 indicates an instance of max pooling utilizing a kernel of  $2 \times 2$  with a stride of 2.



**Fig. 17** Average Pooling using a  $2 \times 2$  kernel with a stride of 2

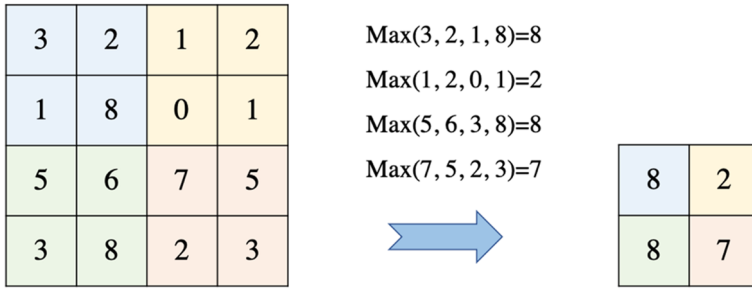


Fig. 18 Max Pooling using a  $2 \times 2$  kernel with a stride of 2

pooling, the maximum value in the  $2 \times 2$  kernel will be left in a rectified feature map. Eventually, the  $4 \times 4$  input image is shrunk to  $2 \times 2$ .

Generally speaking, max pooling is more effective and commonly used in image processing. As it feels more like doing a feature selection. It selects features with better classification and identification as well as providing nonlinearity. Besides, max pooling is able to deal with the second error mentioned in the average pooling paragraph, which means it can preserve more texture information.

For instance, if the feature we need to find is a dog, then as long as there is a dog in one area of the image, it means that there is a dog in the whole image. Therefore, we should take the maximum matching degree between all areas in the image and the characteristics of the dog, that is, maximum pooling should be utilized.

### 3.2.3 Stochastic pooling

In stochastic pooling, the pool mapping response is obtained by sampling the multinomial distribution formed by the activation of each pool area. To be more exact, we compute the probabilities  $p$  for each region  $j$  via normalizing the activations within the region using:

$$p_i = \frac{a_i}{\sum_{k \in R_j} a_k} \tag{31}$$

Where  $a$  refers to value of element.

$k$  refers to the index of element in region of feature map.

We then sample from the multinomial distribution according to  $p$  to choose a location within the region. The procedure is in Fig. 19.

The samples for each pooling region in each layer for each training example are drawn independently to one another. When back propagation is carried out over the network, the

SUM(3, 0, 1, 1, 0, 0, 0, 5, 0)=10    Divide each by 10 (SUM)    Sample a location from P()

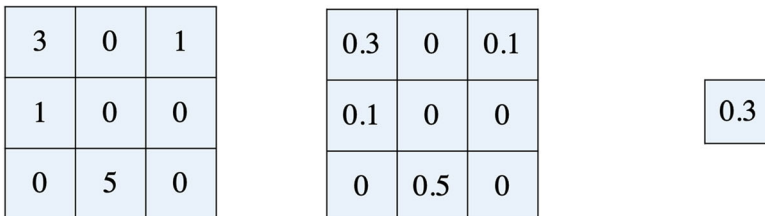


Fig. 19 The transformation process of stochastic pooling

gradient is guided back to the pool area using the same selected location, similar to using the maximum pool for back propagation. When passing information through the network, max pooling only focuses on the strongest activation of each region. While, sometimes non-maximal activations should also be taken into consideration. Stochastic pooling does well in this as those additional activations could be taken into account.

In a word, stochastic pooling is between the average pooling and max pooling. Probability is assigned to pixel points according to the value size, and then sub-sampling is conducted based on the probability. In the average meaning, it obeys average pooling criterion, and in the local meaning, it is similar to max pooling [104].

Although pooling operation has little effect on improving the overall accuracy, it plays an obvious role in decreasing parameters, reducing overfitting, optimizing model performance and saving computational force, so it is an essential process in neural network design [101].

### 3.3 Fully connected layer

Fully connected (FC) layers often appear near the end of CNN system and are utilized to sum the features of the previous design. If we consider convolution and pooling in the front as process of feature engineering, local amplification and extraction of local features. The FC layer in the back can be regarded as feature weighting. A sample structure of FC layer as displayed in Fig. 20 is usually a fast method for learning nonlinear combinations of the high-level properties generated from convolutional layer. In that space, the FC layer will learn about a possible nonlinear function.

The basic procedure of learning is as follow. Firstly, flatten the image that has already been converted to a form suitable for a multilevel perceptron into a column vector and feeding it back to the feedforward neural network. Then, apply the propagation after planarization to each iteration of training. In this way, the model will have the ability to distinguish between the main features and some low-level features in the image and classify them through classification techniques such as SoftMax [31].

In some circumstances, FC layer can be realized by conducting convolution operation, that is, we can utilize a convolutional layer with  $1 \times 1$  kernel size to replace the FC layer that is fully connected to the front layer. Within normal FC layer, the important feature weight is

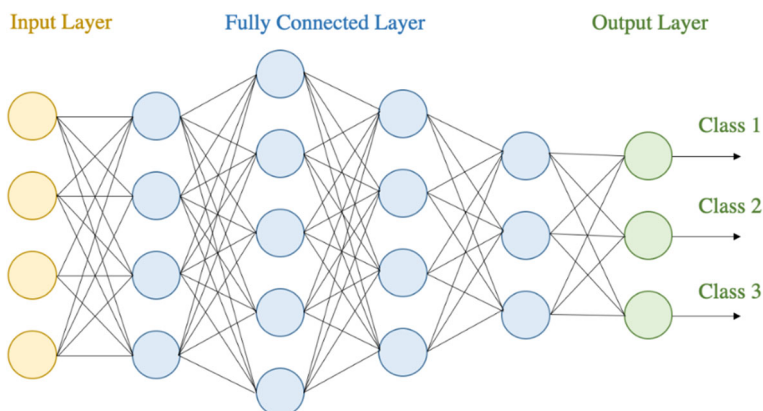


Fig. 20 A sample structure of FC layer

increased to obtain the category of the identified object. While in convolution operation, like an intentional weakening of fully connection, the weak influence outside the local area is directly erased to zero on the basis of the inspiration of local vision. There is also a bit of enforcement that the parameters utilized in different parts are the same. The weakening reduces the parameters, saves the calculation, and specializes in local areas rather than seek for more. While the enforcement can help further reduce the parameters.

After introducing the basic structure of CNN, the next section will focus on some classic CNN architectures.

## 4 Classic architecture

In the early stage, because of the lack of training data and computational ability, it is impossible to train high performance CNN under no over-fitting circumstance. However, in the past few years, with the introduction of ReLU and dropout, along with the historical opportunity achieved by GPU and big data, the research on CNN has experienced several historic breakthroughs, began to obtain first-class results and attracted attention of the world. Figure 21 displays several classic architectures of CNN in its development history [103].

### 4.1 AlexNet

AlexNet came from a team led by Krizhevsky, Sutskever [48] in 2012 and its brief structure is illustrated in Fig. 22. Due to the limited GPU capacity at that time, dual-GPU training was used, and the network was cut into two parts accordingly. However, with the development of GPU, single GPU can easily train AlexNet, so the two parts can be combined as shown in the figure. (Note that this leads to more connections and more computation) From this figure, it is obvious that the entire network structure of AlexNet consists of 5 convolutional layers and 3 FC layers, with a total depth of 8 layers [34, 102].

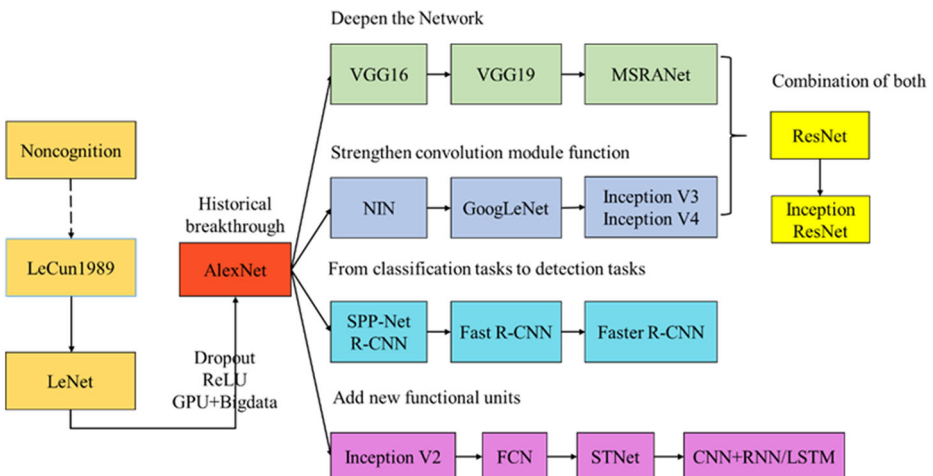


Fig. 21 Classic CNN architecture in development history

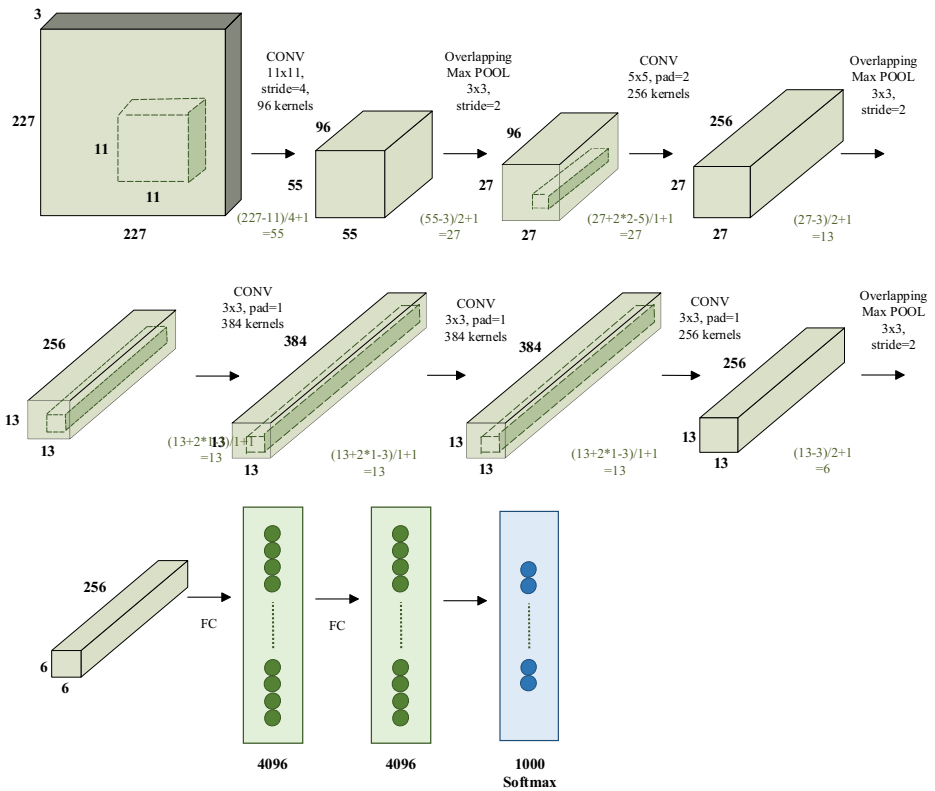


Fig. 22 A sample structure for AlexNet

There is a clear process of direct conversion between layers and a simple interpretation of each layer is displayed in Fig. 23. And then during a dropout process, where you randomly drop some information, you will get 4096 new neurons. The use of dropout can reduce overfitting and enhance the generalization of the model because it makes the model depend less on some local features. The dropout rate usually equals to 0.5.

### 4.2 VGG-net

VGG-Net is a well-known deep CNN developed by Simonyan and Zisserman [89]. It writes about the relationship between the depth and performance of CNN [60]. Through repeatedly stacking the convolutional layer and the maximum pooling layer, VGG-Net creates a 16 ~ 19 layers CNN which achieves the runner-up in ILSVRC 2014 and the champion in positioning, with an error rate of 7.5% on the top5. So far, VGG-Net has still been a widely used network to extract image features.

The innovation of VGG-Net is it uses 3 × 3 convolution kernel and 2 × 2 pooling kernel to improve performance through deepening network structure. The growing number of network layers does not result in an explosion in the number of parameters due to the reason that the number of parameters is concentrated in the final three FC layers. Meanwhile, the series of two 3 × 3 convolutional layers can be regarded as one 5 × 5 convolutional layer, and the series of three 3 × 3 convolutional layers can be regarded as one 7 × 7 convolutional layer, that is, the

<b>CONV1</b>		
Input image specification: 224x224x3 (RGB image), will actually be changed to 227x227x3 after preprocessing.		
96 filters with size and specification of 11x11x3 are used for feature extraction. Use ReLU to ensure that the value of feature graph is within a reasonable range. Data after convolution is 55x55x96.		
Down sampling operation pool1 (Pool1 kernel: 3x3 step size: 2) Data is reduced to 27x27x96.		
<b>CONV2</b>		
Input data 27x27x96.		
256 5x5 filters with step size 1 were utilized to further extract features from the 27x27x96 feature graphs.		
Down sampling operation pool2 (Pool1 kernel: 3X3 step size: 2)		
<b>CONV3</b>		
Get 384 new feature maps of 13x13 (kernel 3x3, step size 1)		
<b>CONV4</b>		
Same as CONV3.		
<b>CONV5</b>		
Output data 13x13x256.		
Pool3 for down sampling operation (Pool3 kernel: 3x3 step size: 2) Data reduced to 6x6x256 [(13-3)/2+1=6]		
<b>FC6</b>	A full connection was made with 4096 neurons to 256 feature maps of size 6 by 6.	
<b>FC7</b>	Similar to FC6.	
<b>FC8</b>	A thousand neurons were used to fully connect the 4,096 neurons in FC7, which then went through a gaussian filter to get 1,000 float values, which is what we consider as predictive probability.	

Fig. 23 Interpretation of each layer in AlexNet

size of the field of three  $3 \times 3$  convolutional layers is equivalent to one  $7 \times 7$  convolutional layer. But the number of convolutional parameters of  $3 \times 3$  is only about half of  $7 \times 7$ , and the former can have 3 nonlinear operations while the latter only has 1 nonlinear operation, which makes the former have relatively stronger learning ability for features than the latter.

During training, VGG-Net first trains network of level A, and then initializes the later complex models by repeatedly using the weight of network A, in this way, the convergence speed improves. The configuration of VGG-Net is shown in Fig. 24.

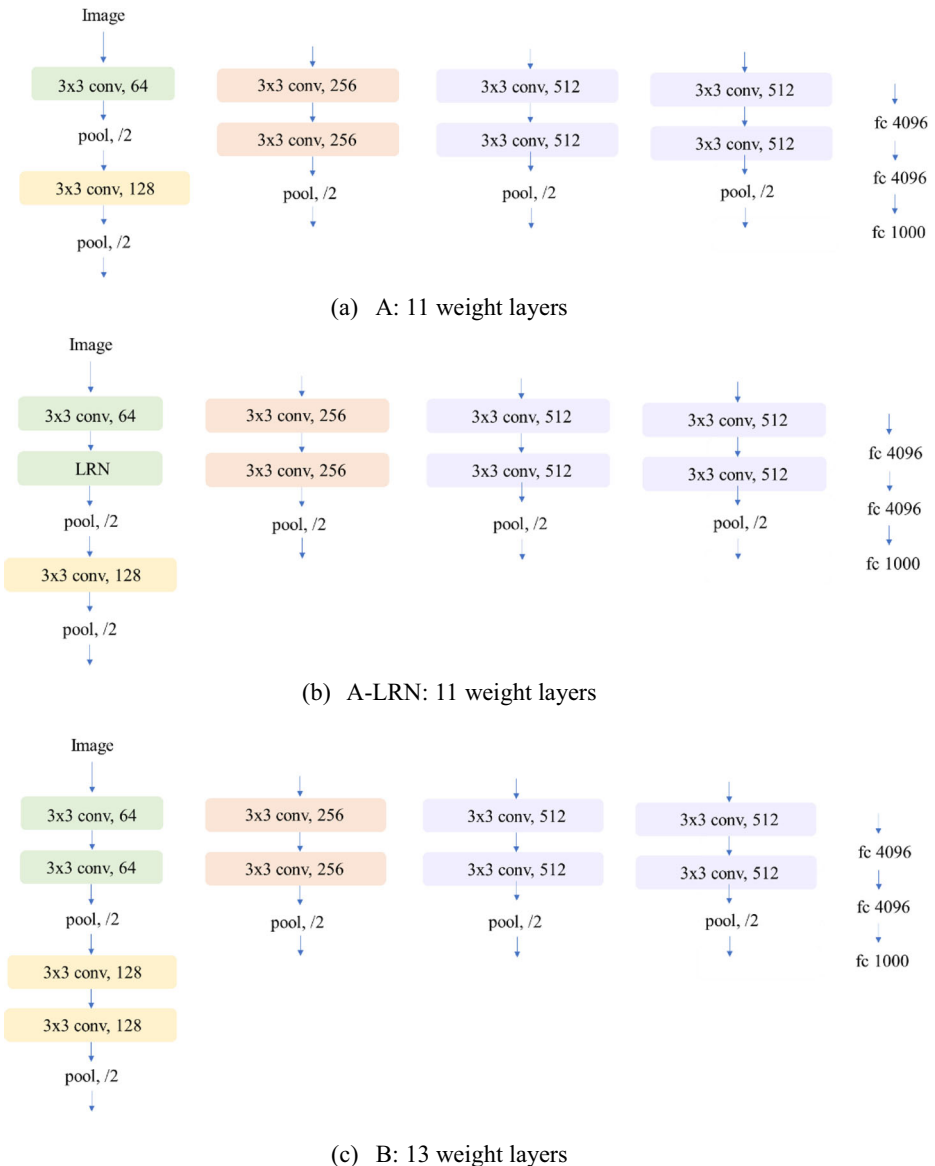
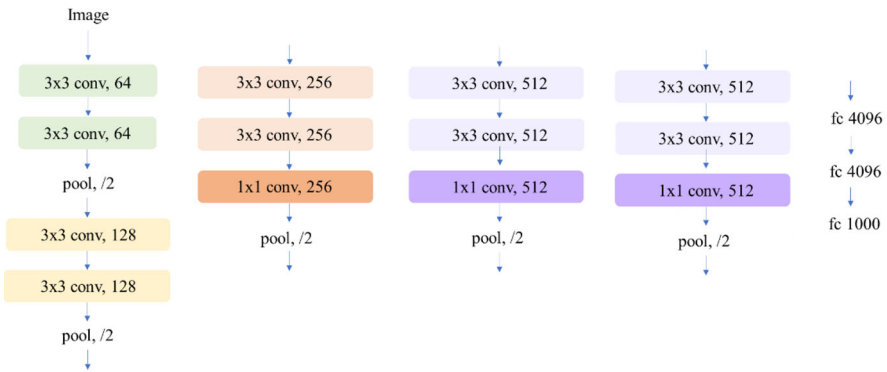
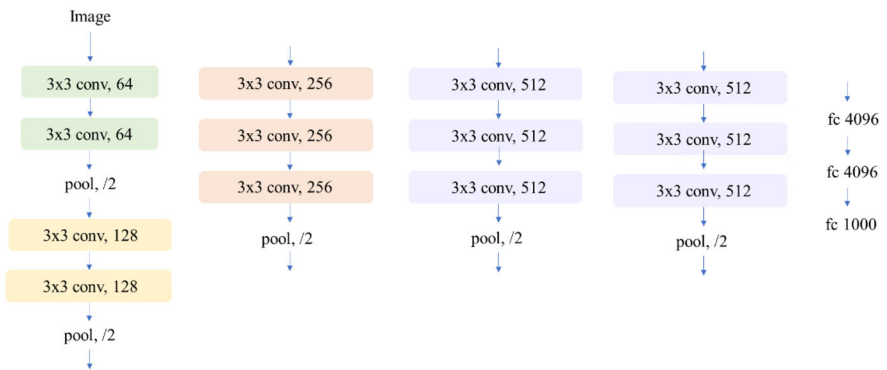


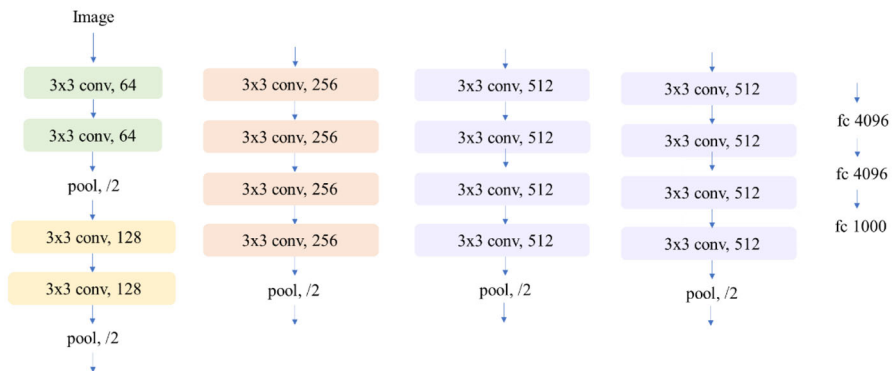
Fig. 24 Configuration of VGG-Net



(d) C: 16 weight layers



(e) D: 16 weight layers (VGG-16)



(f) E: 19 weight layers (VGG-19)

Fig. 24 (continued)

The authors of VGG-Net conclude that LRN layer has little effect. The experiment confirms that, when the network is going deeper, performance of network would be improved



as well. The convolution of 1 by 1 is also efficient, allowing linear transformations of pixels without changing the number of output channels and the parameters of the neural network. But it doesn't work as well as the convolution of 3 by 3, because the convolution of 3 by 3 learns more about spatial features.

### 4.3 Inception (GoogLeNet)

Inception (also known as GoogLeNet) is a classic deep learning structure proposed by Szegedy, Ioffe [95]. Before this, many structures mentioned before all utilize adding the depth of the network to achieve better training performance, however, utilizing more layers will at the same time bring loads of negative effects, for instance, overfitting, gradient disappearance, gradient explosion, etc. Inception improves the training results by proposing a new way which is making more efficient use of computing resources, in other words, extracting more features for the same amount of computation [93, 94].

The basic structure is illustrated in Fig. 25. The entire inception structure is made up of several such inception modules in series. The inception structure has two main contributions: one is using  $1 \times 1$  convolution to raise and lower the dimension, the other is carrying out convolution and reintegration on multiple scales at the same time.

There are multiple yellow  $1 \times 1$  convolution blocks in the figure, which have two effects.

The first effect is they can help extracting richer features by adding more convolution to the same size of the feeling field. This idea comes from the Network in the Network (NIN) which indicates that if you put two convolutions in series, you can combine more nonlinear features. The structure of NIN is likely to the multi-layer structure in the traditional neural network. The latter multi-layer spans the receptive field of different sizes by adding pooling layer between layers, so as to extract features at a higher scale. NIN structure is multi-layer in the same scale (without pooling layer in the middle), so it can extract stronger nonlinearity in the same receptive field [57].

The second effect is  $1 \times 1$  convolution can be used to reduce the dimension, which helps decreasing the computational complexity. More specifically, when a convolution layer has loads of input features, the convolution operation on the input will produce much computation. However, if shrinking the dimensionality of input first, and doing the convolution after the number of features are reduced, the computation will significantly be reduced as well. Figure 26 compares the multiplication times of the two schemes before and after optimization.

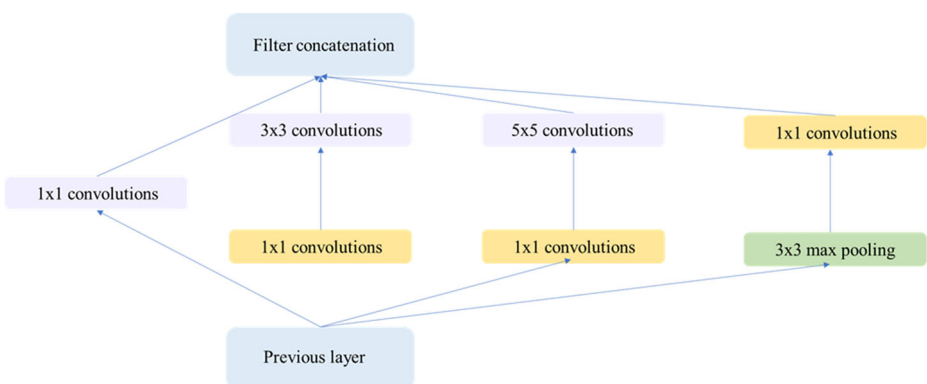
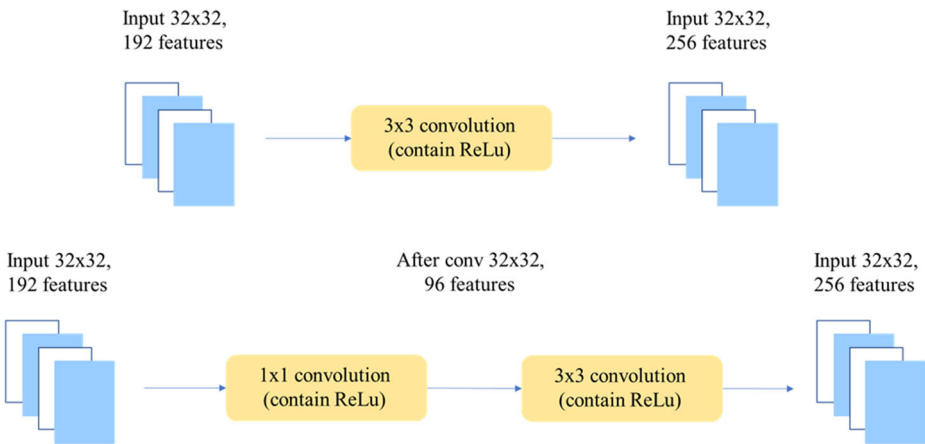


Fig. 25 Basic structure of the inception module



**Fig. 26** The comparison between the multiplication times of the two schemes before and after optimization

With the same input of a set of 192 features,  $32 \times 32$  size, and the output of 256 features, the first scheme is directly realized by  $3 \times 3$  convolution, which requires.

$$192 \times 256 \times 3 \times 3 \times 32 \times 32 = 452,984,832 \text{ multiplications} \quad (32)$$

In the second scheme, the convolution of  $1 \times 1$  was first utilized to decrease the number of features to 96, and then  $3 \times 3$  was used to restore 256 features, which required.

$$\begin{aligned} 192 \times 96 \times 1 \times 1 \times 32 \times 32 + 96 \times 256 \times 3 \times 32 \times 32 \\ = 245,366,784 \text{ multiplications} \end{aligned} \quad (33)$$

The dimensionality reduction method of  $1 \times 1$  was used to save half of the computation. Besides, if the number of features in the final output remains the same (256 sets), the dimensionality reduction in the middle is similar to the effect of compression and will not influence the result of the eventual training.

Apart from the  $1 \times 1$  convolution blocks, the input layer has four branches, which were convolved or pooled with filters of different sizes, and finally integrated together on the feature dimension. Szegedy explained the benefits of this new structure from a number of perspectives.

The first explanation is the features of different scales can be extracted by convolution on multiple scales at the same time. The richer the features, the more accurate the final classification.

The second explanation is utilizing the principle of decomposition of sparse matrix into dense matrix to accelerate the convergence. For instance, the first matrix in Eq. (34) is a sparse matrix (many elements are 0 and are unevenly distributed in the matrix), which is convolved with a  $2 \times 2$  matrix and needs to be calculated for each element in the sparse matrix. However, if the sparse matrix is decomposed into 2 sub dense matrices and then convolved with the  $2 \times 2$  matrix, the regions with more 0 in the sparse matrix can be eliminated from calculation, greatly reducing the amount of calculation. Applying this principle to Inception needs decomposing on the feature dimension. The input data of the traditional convolution layer convolved with the convolution kernel of one scale (such as  $3 \times 3$ ) and output fixed dimensions data (such as

256 features). All 256 output features were uniformly distributed on the scale of  $3 \times 3$ , which can be understood as the output of a sparsely distributed feature set. While Inception module extract features on multiple scales (such as  $1 \times 1, 3 \times 3, 5 \times 5$ ), thus the output features of 256 is not evenly distributed, but the features that have a strong correlation gather together (such as  $96 \ 1 \times 1$  features together,  $96 \ 3 \times 3$  features together and  $64 \ 5 \times 5$  features together). This can be understood as multiple densely distributed sub-feature sets. In this feature set, as the more relevant features cluster together, the non-relevant non-key features weaken, so that the Inception method outputs the same 256 features with less ‘redundant’ information. With such ‘pure’ feature set transferred layer by layer and eventually used as input of reverse calculation, the natural convergence rate becomes faster.

$$\begin{pmatrix} 3 & 4 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 3 & 0 & 0 \\ 0 & 0 & 0 & 1 & 4 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \otimes \begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix} \Leftrightarrow \frac{\begin{bmatrix} 3 & 4 \\ 0 & 2 \end{bmatrix} \otimes \begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix}}{\begin{bmatrix} 2 & 3 & 0 \\ 0 & 1 & 4 \\ 1 & 0 & 0 \end{bmatrix} \otimes \begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix}} \tag{34}$$

Where  $\otimes$  represents Kronecker Product which is used for operations between two matrices of any size.  $\Leftrightarrow$  represents Substantively equivalent.

The third explanation is according to Hebbin’s principle [81]. The Hebbin’s principle is a theory in neuroscience that explains what happens to neurons in the brain during the learning process: fire together, wire together. Hebbin believes that ‘Two neurons, or systems of neurons, that are always firing at the same time form a ‘combination’ in which the firing of one neuron promotes the firing of the other’. For example, dogs drool when they see meat. After repeated stimulation, the neurons in the brain that recognize meat and those responsible for salivating will promote each other and ‘entwine’ with each other, making them drool faster when they see meat in the future. The usage of this principle in Inception structures is to bring together highly relevant features. It’s a little bit like the second explanation above, where you separate the features of  $1 \times 1, 3 \times 3$ , and  $5 \times 5$ . Since the ultimate goal of training convergence is to extract independent features, the convergence speed can be accelerated as long as the highly relevant features be gathered in advance.

### 4.4 ResNet

ResNet was proposed by He, Zhang [27] from Microsoft research institute in Dec, 2015. Through utilizing residual Unit, ResNet successfully trained the 152-layer neural network, and achieved the best result in ILSVRC 2015. Meanwhile, the number of parameters was lower than VGG-Net, with outstanding effect. The structure of ResNet is able to greatly speed up the training of ultra-DNN, as well as perfect the model accuracy.

The initial inspiration of ResNet comes from this problem: if depth of neural network is continuously increasing, a degradation problem will occur. This means the accuracy will first rise and then reach saturation, afterwards the continuous increase of depth will cause the decline of accuracy. The difference between this problem and overfitting is, in this problem, the error increases in both testing set and training set.

If a relatively shallow network reaches saturated accuracy, then adding several  $y=x$  congruent mapping layers to the end will at least not increase the error, which means, the

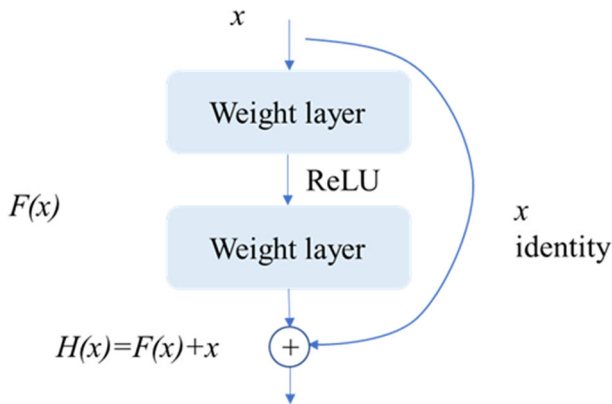


Fig. 27 A sample structure of residual block

deeper network will not cause the increasing of error on training set. The idea of utilizing congruent maps to transfer the previous layer of output directly to next layer was the inspiration for ResNet. If we directly pass the input to the output as the beginning value, the target we require to study at this point will be:

$$F(x) = H(x) - x \quad (35)$$

Where  $x$  refers to input of neural network.

$H(x)$  refers to expected output

The Residual Unit of ResNet is illustrated in Fig. 27. ResNet altered the learning objective from learning a complete output to learning only the difference between output and input, thus simplifying the learning objectives and difficulty.

The Residual block is realized by shortcut connection as displayed in Fig. 27. Through utilizing shortcut, this block performs an element-wise superposition of the input and output. This simple addition will not increase additional parameters and computation to network. However, it could efficiently improve the training speed and training effect of the model. In addition, if the layer number of model is increasing, this structure is also able to efficiently handle in the problem of degradation.

Figure 28 shows a comparison between a 34-layer ordinary CNN and a 34-layer ResNet network. It can be seen clearly that the biggest difference between a straight-through CNN and ResNet is that ResNet has a large amount of bypass branches which transmit the input straightway to the subsequent layer, allowing the subsequent layer to directly study the residual. This category of structure is as well-known as shortcut or skip connection.

## 4.5 U-net

U-Net is a classical full convolutional network. This U-shaped structure can capture semantic information in the process of down-sampling and accurately locate the corresponding up-sampling process. This network can train end-to-end with a small number of images, but the effect exceeds that of the previous best model, CNN based on sliding Windows [74]. The U-shaped structure is shown in Fig. 29.

Assuming that the input image tile is of  $572 \times 572$ . The left side of the network composed of four blocks is called contracting path, which is a series of reduction

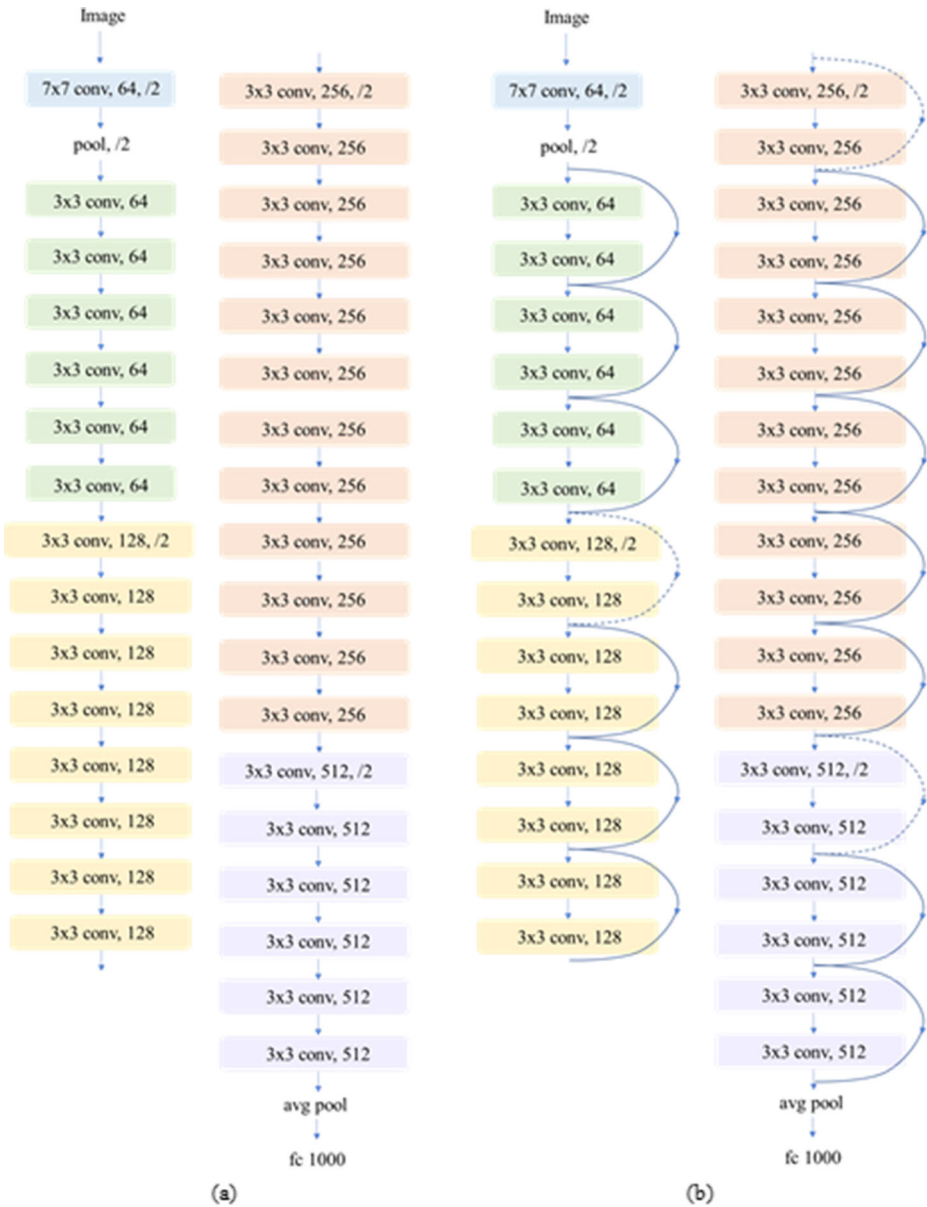


Fig. 28 Comparison between a 34-layer ordinary CNN and a 34-layer ResNet network

sampling operations formed by convolution and max pooling. Finally, a feature map with a size of  $32 \times 32$  is obtained.

Network on the right-hand side is named expansive path. It is as well composed of four blocks. Before the start of each block, it will use deconvolution to enlarge two times the size of the feature map, and the number in half. Then, it is merged with the feature map of the symmetric contracting path on the left due to the reason that the size of the feature map of the contracting path on the left is different from that of the expansive path on the right. The normalization of U-Net is conducted

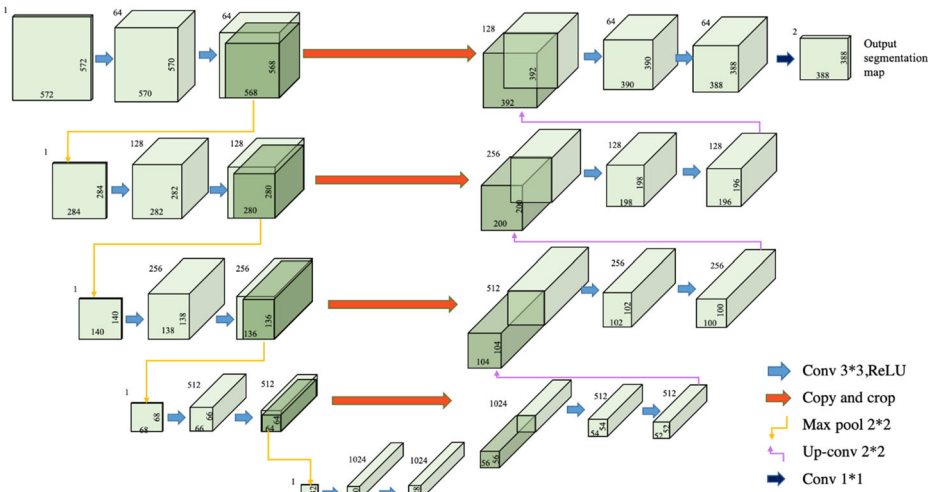


Fig. 29 The basic structure of U-Net

by clipping feature map of the contracting path to the same size as feature map of the expansive path (i.e., the dark green block in Fig. 29). Effective convolution operation is still used in convolution operation of extension path, and the final size of feature map is  $388 \times 388$ .

U-Net mainly uses SoftMax cross entropy loss to deal with the problem that similar target boundaries of medical images are not easy to be recognized. Besides, weights are added to each pixel in the objective function to make the network more able to distinguish the boundaries. The equations are as follows:

$$E = \sum_{x \in \Omega} w(x) \log(p_{l(x)}(x)) \tag{36}$$

$$w(x) = w_c(x) + w_0 \cdot \exp\left(-\frac{(d_1(x) + d_2(x))^2}{2\sigma^2}\right) \tag{37}$$

Where  $p_{l(x)}$  is loss function of SoftMax.

$l$  refers to the label value of a pixel.

$w$  refers to the weight value of a pixel.

$w_c$  is the weight that balances the proportion of the category.

$d_1$  is the distance from the pixel to the nearest cell.

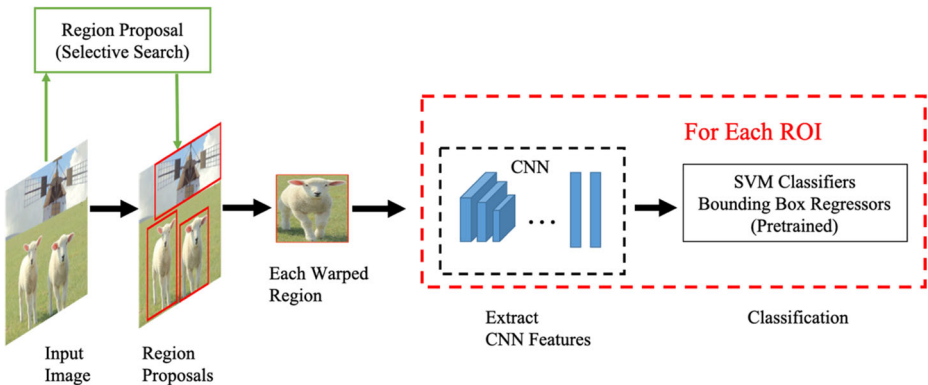
$d_2$  is the distance from the pixel to the second closest cell.

$w_0$  and  $\sigma$  are constants.

## 4.6 R-CNN series

### 4.6.1 R-CNN

The general procedure for target detection of R-CNN is illustrated in Fig. 30: (i) Extract regional proposals from input images using selective search algorithm and warp them to a fixed size. Selective search synthesizes the methods of both continues search and



**Fig. 30** The basic structure of R-CNN

segmentation. It aims at identifying possible target locations for object recognition. (ii) Input the normalized regional proposals into the CNN network and extract features. (iii) SVM classification is used for identification of these features, and linear regression is applied for fine-tuning the border position and size, with each category separately training a bounding box regression. By using less and higher quality region of interest (ROI), R-CNN performs better than the system using traditional sliding window approach in the evaluation of accuracy and speed [24].

However R-CNN framework still exist the following problems: (i) it is very time-consuming and has lots of repeated computation as each candidate box needs to go through CNN alone; (ii) its training is divided into several stages and the steps are tedious; (iii) it takes up disk space: the convolved characteristics of the data needs to be kept separately [53].

#### 4.6.2 SPP-net

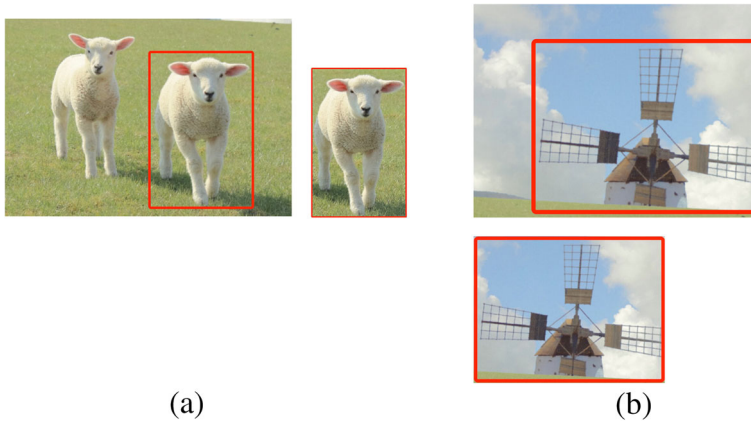
In order to solve the problems of R-CNN, SPP-Net made the following improvements.

One of the main ideas of SPP-Net is to remove the operations such as crop/ warp on the original image and replace them with the spatial pyramid pooled layer (SPP) on the convolution feature. In R-CNN, the proposal was warped before it was sent to the net. Part of the reason for the introduction of SPP layer is that different levels of crop/ warp on images in R-CNN will not only take time and effort but also cause some problems as shown in Fig. 31. For example, the crop in the figure will cause the object to be incomplete, and warp will cause serious deformation after the object is stretched [70].

However, the more important reason is that in R-CNN, each regional proposal is subject to a CNN feature extraction operation, which brings a lot of double calculations. In SPP-Net, we only performed CNN feature extraction on the whole image once, and then map the proposal to the feature map generated by the convolutional network, which is equivalent to feature map obtained by image in a proposal going through the layer by layer convolution process.

#### 4.6.3 Fast R-CNN

It could be obviously seen in Fig. 32 that fast R-CNN is also a method that makes some perfections based on R-CNN. Similar to SPP-Net, it only conducts CNN feature extraction for

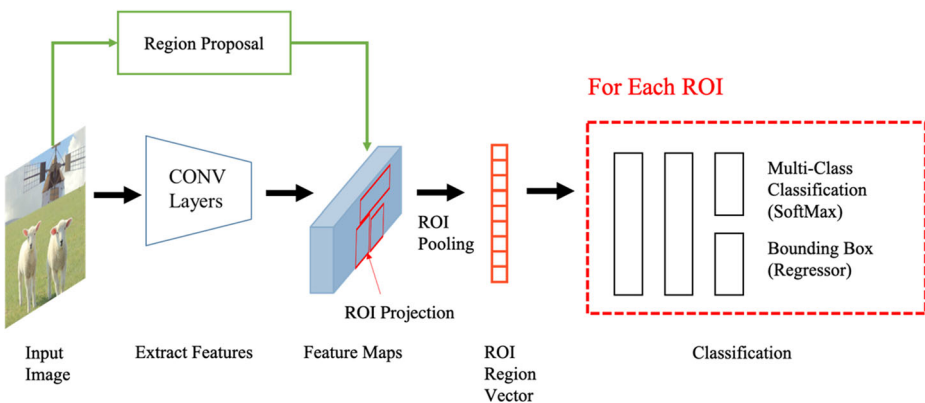


**Fig. 31** (a) Crop; (b) Warp

the whole image once, and then utilizes an ROI pooling layer. While, due to the reason that it is not a fixed size input, the size of the pooling grid shall be adjusted dynamically each time to achieve regional normalization. Fast R-CNN uses SoftMax to replace SVM for classification, as well as multi-task loss function to carry out border regression and classification together [23].

The main steps are as follows: (i) Take the whole picture as the input and use CNN to get the feature layer of the image. (ii) Extract regional candidate boxes from the original image by selective search and projected onto the final feature layer one by one. (iii) ROI pooling operation is conducted for each area candidate box on the feature layer to achieve feature representation of fixed size. ROI pooling is able to accelerate the speed of training and testing. (iv) After two FC layers, utilize SoftMax to conduct multi-target classification respectively, and apply regression model in fine-tuning the border position and size.

The disadvantage of fast R-CNN is that region proposal is extracted by selective search, and most of the target detection time is spent on it (normally, if a region proposal is proposed 2 ~ 3 s once, the feature classification will be proposed only 0.32 s once), which could not reach the requirements of real-time application.



**Fig. 32** The basic structure of fast R-CNN



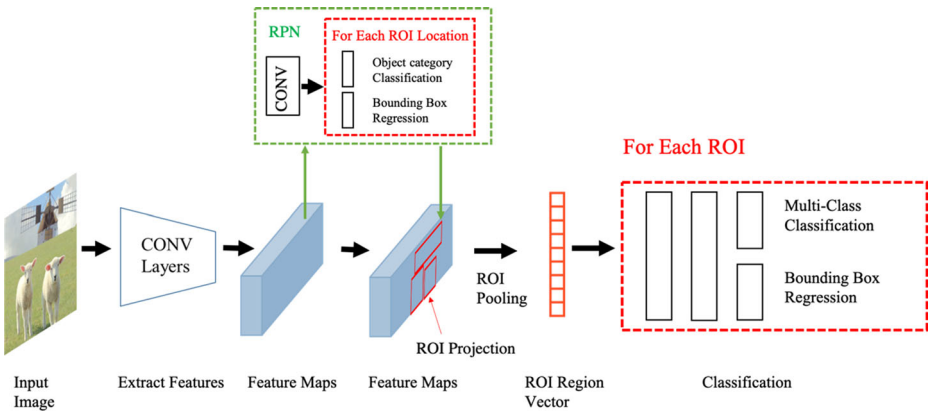


Fig. 33 The basic structure of faster R-CNN

In general, fast R-CNN uses feature CNN as extractors to extract features of the whole image first. And then it applies the method of creating candidate regions directly to the extracted feature map. For example, fast R-CNN chose feature map output by convolutive layer Conv5 in VGG-16 to generate ROI. These ROI would then be combined with corresponding feature graphs to cut out feature blocks and used in target detection tasks. We use ROI pooling to convert feature blocks to fixed sizes and feed them to the FC layer for classification and positioning. Because fast R-CNN has no repeated feature extraction

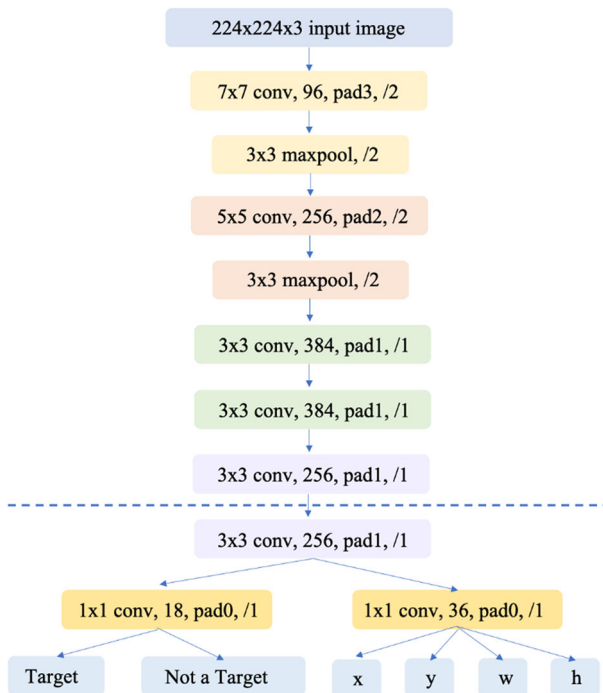


Fig. 34 The basic structure of RPN Network

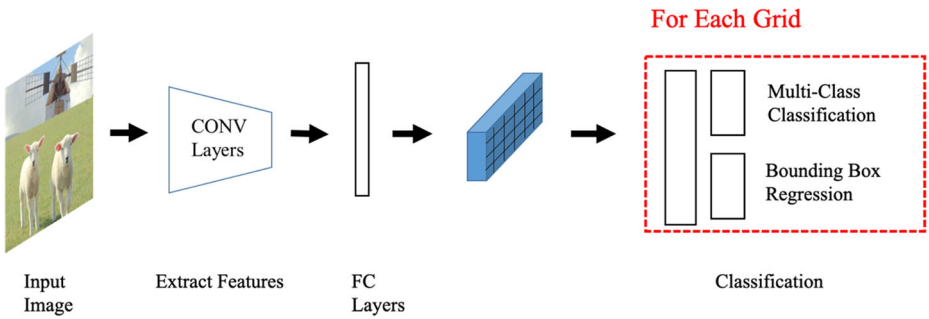


Fig. 35 The basic structure of YOLO

operation, it saves lots of processing time. The training speed of fast R-CNN is 10 times faster than R-CNN, and the predicting speed is 150 times that of R-CNN.

### 4.6.4 Faster R-CNN

As displayed in Fig. 33, faster R-CNN utilizes the proposal achieved by Region Proposal Network (RPN) to replace the proposal achieved by selective search.

The function of RPN network is outputting a batch of rectangular candidate regions with the input of an image, which is somewhat alike to the selective search in the previous target detection. The network structure is according to a CNN, but includes the SoftMax and Bounding box regression multi-task models [73].

In Fig. 34, above the dashed line is the structure in front of the last convolutional layer of ZF network. Under the dashed line is the particular structure of RPN network. For the particular structure, the convolution is 3\*3, and then the output is divided into two ways

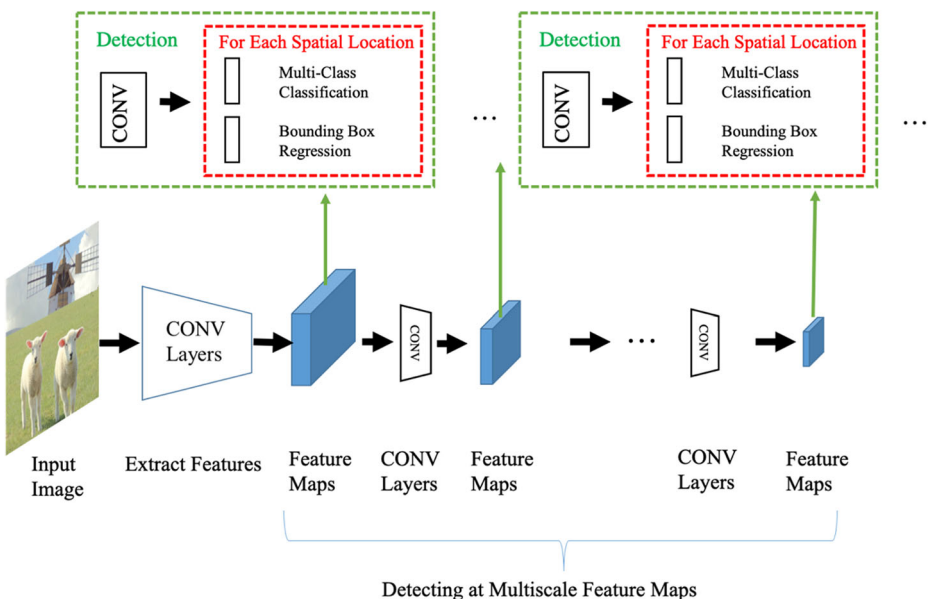


Fig. 36 The basic structure of SSD

through  $1 \times 1$  convolution, one way of which outputs the probability of target and non-target, and the other way outputs the four parameters related to box.

The procedure of faster R-CNN is: (i) Same as fast R-CNN. (ii)  $K$  different Anchor boxes are used in the final convolution feature layer to conduct a region proposal, that is, object/non-object binary classification is carried out for the corresponding area of each anchor box, and  $k$  regression models are used for fine-tuning the location and size of the candidate box. Regional proposal (via RPN) is generated on the feature map, with about 300 suggestion windows per image. (iii) Same as fast R-CNN. (iv) Classify the objects and do border regression.

## 4.7 Yolo

### 4.7.1 YOLOv1

YOLO series algorithm is a typical one-stage target detection algorithm. Its basic structure is illustrated in Fig. 35. YOLOv1 is an end-to-end target detection approach, and the determination of position and category is accomplished by direct regression without need for intermediate regional proposal to find the target.

If two small targets fall into a cell at the same time, the model can only predict one. Large object intersection over union (IOU) error and small object IOU error are close to the loss contribution value in network training. Thus, for small object, small IOU error will also bring a great influence on the network optimization, reducing positioning accuracy of object detection. YOLO uses multiple lower sampling layers, and the object features obtained from network are not precise, which will also affect the detection effect [19].

Compared with faster R-CNN, YOLOv1 has more obvious advantages. Firstly, faster R-CNN needs an RPN network to find the candidate region, while for YOLOv1, after the input image has been divided, it directly takes the region as the candidate region. YOLOv1 utilizes the end-to-end regression method rather than use those regional proposal steps, so direct regression accomplishes the determination of position and category. Therefore, compared with fast R-CNN, the main advantage of YOLOv1 lies in (i) its speed; (ii) global processing makes background errors less than regional based methods; (iii) good generalization performance. However, YOLOv1 performs not very well in target positioning process, which causes low detection accuracy of YOLO.

To sum up, YOLOv1 proposes a new idea of a single stage. Compared with the two-stage method, YOLOv1 has obvious advantages in speed and impressive real-time characteristics. However, YOLOv1 also has some problems, such as rough grid division and the number of boxes that could be generated by each grid, which limits the detection of small-scale objects and objects that are close to each other.

### 4.7.2 YOLOv2

The differences between YOLOv2 and YOLOv1 are as follows:

- 1) The dropout layer was removed, and Batch Normalization was added in each convolutive layer, with a 2% improvement in Mean Average Precision (mAP).
- 2) High-resolution classifier. Input data of high-resolution changed from  $224 \times 224$  to  $448 \times 448$ , and mAP improved by 4%.

- 3) YOLOv2 deletes a pooling layer to get high resolution data. Use anchor to replace FC layer.
- 4) Anchor box is determined by k-mean clustering and distance calculation of clustering is used instead of Euclidean distance.
- 5) YOLOv2 uses anchor box to predict the object category rather than use grid to predict the probability of each grid in YOLOv1.
- 6) The bounding Box coordinate continues YOLOv1's prediction of Bounding boxes' coordinates with respect to the upper left corner of the grid. To ensure that the center of bounding Box falls within the grid (i.e. its coordinates fall between 0 and 1), Bounding Box center is sent into the activation function.
- 7) YOLOv2 adds a passthrough layer to improve the detection of small objects. This layer conducts information fusion of high-resolution feature and low-resolution feature. Through the stacked in adjacent features rather than the spatial location of different channel, it integrates characteristic of high and low resolution, making the mAP value increased by 1%.
- 8) YOLOv2 only uses convolutional layer and pooling layer, thus the size of input image could be dynamically adjusted. Besides, the accuracy and speed can be adjusted according to own requirements.

## 4.8 SSD

YOLO predicts a location utilizing full-image features, SSD predicts a location utilizing features around that location [58].

According to Fig. 36, it seems that SSD model is somewhat alike to RPN, for instance, the Default Bounding Box in SSD is very likely to be the anchor in RPN. However, SSD takes into account different scales in different feature layers, and RPN thinks about different scales in one feature layer.

SSD outputs a series of discretization boxes, produced on feature maps on different layers, with different aspect ratio.

This section introduced the classic architectures in CNN from an early stage in nearly 10 years ago till now. The CNN model has been continuously evolved and eliminated in the careful study of countless scientists, thus forming various model systems with their own merits. The following section will focus on how different model made remarkable achievements and outstanding contributions in different medical fields.

## 5 Application

Since CNN have some many models and it can be used in lots of fields of image classification like identifying human face and patient condition. The following content will introduce the application of CNN for classifying the patient conditions in different body parts.

### 5.1 Breast

According to the statistical analysis defined from Institute for Health Metrics and Evaluation (IHME), there was a noticeable increase in breast cancer mortality in most regions of the world

between 1990 and 2015, with a total mortality rate of 0.7 per 100,000. Data also indicates that the five-year survival rate of early stage breast cancer patients exceeds 95%, while the five-year survival rate of late period breast cancer patients is less than 20%. The World Health Organization (WHO) has confirmed that early diagnosis and treatment of breast cancer patients could not only decrease mortality, but also help some patients avoid psychological suffering and financial loss from mastectomy and axillary lymph node dissection [5, 83, 98]. Thus, early detection is a critical step in prevention and treatment of breast cancer.

Mammography is a specific imaging method to screen and diagnose breast cancer. Lévy and Jain [54] and his group used CNN to straightway classify pre-detected breast masses in mammograms as benign or malignant. They utilized AlexNet, GoogLeNet model and achieve outstanding precision rate result as shown in Table 1. Another project focuses on applying CNN to detect microcalcifications in mammograms, since CNN shows excellent performance in classifying image data. In order to solve the problem of large class imbalance between pixels from microcalcifications and other breast tissue, the research team of Mordang, Janssen [66] applied a CNN architecture which is inspired by OxfordNet to classify the microcalcification pixels which are difficult to be identified. The researchers also used cascade classifier for comparison and find out that CNN classifier had significantly better results for different data set with a sensitivity increase of about 6%. Except the Mammography, breast cancer cytological specimen can also be analyzed in order to diagnose benign and malignant cases. To improve the diagnosis of breast cancer, team of Gao, Wu [21] presented a network in combination of a shallow CNN and a deep CNN. They developed the shallow one to achieve ‘virtual’ recombined images from input dataset and apply the deep CNN to extract features. Finally, the recombined and ‘virtual’ recombined images would help the model to classify benign or cancer category of brain circumstance. In the research of Sun, Tseng [92], the group applied CNN into a semi-supervised learning scheme for breast cancer detection, which only needs very few labelled data for training. This scheme allows the usage of unlabelled data if a

**Table 1** Results of application of convolutional neural network in breast area

Authors	Goals	Model	Results	
			Acc	Sen
D Lévy et al. [54]	Classify the mass, benign or malignant	AlexNet GoogLeNet	0.8900 0.9290	0.8680 0.9340
J Mordang et al. [66]	Detect microcalcifications and then classify the mass, benign or malignant	Two CNNs	/	0.7035 ±0.0068
F Gao et al. [21]	Classify the benign or cancer	Shallow-Deep CNN	0.90 ± 0.06	/
W Sun et al. [92]	Allow the usage of unlabelled data in breast cancer diagnosis	A CNN model including 3 pairs of convolutional layers	0.8818	/
H Su et al. [90]	Relieve computational burden in segmentation by reducing the redundant computations	Fast scanning CNN	0.9100	/
F Ting et al. [96]	Classify incoming medical images as malignant, benign or normal patient	CNNI-BCC	0.9050	0.8947
D Bardou et al. [7]	Binary classification Multi-class classification	CNN using an ensemble model	0.9833 0.8233	/

(Acc = Accuracy; Sen = Sensitivity)

small amount of labelled data is already learned by the system. The final result achieves 88.18% with unlabelled data and 82.43% with the mixed and unlabelled data, which shows a significant improvement than previous research in accuracy. In the pixel-wise region segmentation scheme proposed by Su, Liu [90], fast scanning deep CNN was applied to relieve the computational burden by reducing the redundant computations in the convolutional layers and max pooling layers. Ting, Tan [96] introduced a model for breast cancer diagnosis and named it as Convolutional Neural Network Improvement for Breast Cancer Classification (CNNI-BCC). This algorithm achieves 89.47% sensitivity and 90.50% accuracy. Bardou, Zhang [7] tried to figure out the breast cancer classification on the basis of histology images utilizing CNN. The result reveals that CNN performs better than the handcrafted feature-based classifier with 98.33% for binary classification and 82.33% for multi-class classification.

## 5.2 Lung

Lung cancer is a major global health problem today and for decades to come, as it is the major cause of cancer deaths worldwide, far outnumbering breast cancer, accounting for 18.4% of all cancer deaths. The data show that about 70% of patients are diagnosed with advanced lung cancer, so only 15% of patients with lung cancer survive five years after being diagnosed. This means that early diagnosis and improved diagnostic accuracy are crucial to the treatment of lung cancer and this area of technology has long been the subject of research [12, 18].

Alakwaa, Nassef [2] with his team used CT scans as input and 3D-CNN as model to detect and classify the lung unmarked nodules. They first segmented out lung tissue by thresholding, and then, a modified U-Net that can determine nodule candidates in CT scans is used to locate the regions that has relatively high possibility of containing nodule candidates. Finally, these regions with nodule candidates were sent to 3D-CNN to be identified if there exists lung cancer, since the 3D-CNN could make fully usage of the spatial 3D context information of lung nodules. The accuracy of model achieved 86.6%. In another similar research displayed in Table 3 in detail, the group of Kang, Liu [45] added the concept of multi-view strategy into the usage of 3D-CNN to improve the model performance. They conducted a binary classification and a three-type classification on CT images and compared results among different methods. Compared with 2D

**Table 2** Results of application of convolutional neural network in lung area

Authors	Goals	Model	Results	
			Acc	Sen
W Alakwaa et al. [2]	Detect the lung unmarked nodules	U-Net and 3D CNN	0.867	0.853
G Kang et al. [45]	Binary classification	MVO-3D CNN-CNN1	0.9541	/
	Ternary classification		0.923	/
W Shen et al. [87]	Model raw nodules to detect nodule heterogeneity	Multi-scale CNN	0.8684	/
A Chris et al. [4]	Lung Pattern Classification	Self-created	0.8561	/
H Sahiner et al. [26]	Automatically detect lung nodules	3D FCN + CNN	/	0.95
H Gillies et al. [33]	Lung nodule characterization	TUMORNET	0.9231	/

(MVO = multi-view-one; Acc = Accuracy; Sen = Sensitivity)

\*(CT1 = CT image channel 1; CT3 = CT image channel 3)

CNN, 3D CNN is no doubt better in error rate whatever in binary classification or ternary classification. Moreover, compared with one-view-one-network strategy, the multi-view-one-network strategy also achieves a lower error rate. The final best results achieved a superior error rate of 4.59% for the binary classification and 7.70% for the ternary classification. One breakthrough brought by Shen, Zhou [87] and his team was that, in terms of regional analysis, instead of depending on nodule segmentation, they directly modeled raw nodule patches without any prior definition of nodule morphology. A Multi-scale Convolutional Neural Networks (MCNN) model was set up to capture nodule heterogeneity and the input also comes from CT which shows both lung nodule screening and nodule annotations. The best accuracy rate achieved 86.84% with the RF classifier as displayed in Tables 2,3 and 4. The first model for dealing lung pattern classification problem was presented by Anthimopoulos, Christodoulidis [4] and proven to be effective when compared to previous methods on a public dataset. In the research of Hamidian, Sahiner [26], they first presented a fully convolution network (FCN) transferred from 3D CNN and used it to produce scoring results for all the volumes rather than just volumes of interest. Then, utilizing the difficult negative examples produced by FCN to train a new discriminant CNN model. The whole system was a

**Table 3** Details in research of G Kang et al. [45]

Input data	Goals	Model			Acc	Sen	Spc		
CT 1	binary classification	MVO	3D CNN	CNN1	0.9518	0.9556	0.9301		
				CNN2	0.9503	0.9551	0.9299		
				CNN3	0.9481	0.9845	0.8976		
		2D CNN	CNN1	0.9454	0.8888	0.9967			
			CNN2	0.944	0.8893	0.9972			
			CNN3	0.9413	0.894	0.9874			
			CT 3	3D CNN	CNN1	0.9525	0.956	0.9394	
					CNN2	0.9524	0.9561	0.9389	
					CNN3	0.9362	0.9417	0.8973	
2D CNN	CNN1	0.9459	0.8892	1					
	CNN2	0.9446	0.89	0.9976					
	CNN3	0.944	0.8897	0.9968					
CT		MVO OVO	3D CNN	CNN	0.9541	/	/		
				CNN1	0.9472	/	/		
				CNN2	0.9439	/	/		
CT 1	ternary classification	MVO	3D CNN	CNN1	0.9188	/	/		
				CNN2	0.9158	/	/		
				CNN3	0.8906	/	/		
			2D CNN	CNN1	0.8234	/	/		
				CNN2	0.8373	/	/		
				CNN3	0.8479	/	/		
				CT 3	3D CNN	CNN1	0.9222	/	/
						CNN2	0.9175	/	/
						CNN3	0.9069	/	/
2D CNN	CNN1	0.8184	/	/					
	CNN2	0.8195	/	/					
	CNN3	0.8681	/	/					
CT		MVO OVO	3D CNN	CNN	0.923	/	/		
				CNN1	0.9187	/	/		
				CNN2	0.9093	/	/		

(MVO = multi-view-one; OVO = one-view-one; Acc = Accuracy; Sen = Sensitivity; Spc = Specificity)

\*(CT1 = CT image channel 1; CT3 = CT image channel 3)

**Table 4** Results of application of convolutional neural network in brain area

Authors	Goals	Model	Results			
			Acc	Dss		
				Com	Co	E
K Shah et al. [46]	Fully automatically segment the brain the brain tumor	TP-CNN**	/	0.81	0.76	0.73
S Khan et al. [79]	Classify multi-grade brain tumor	FT-VGG-19	0.9458	/		
A Ayati et al. [3]	Classify different grades of glioma	CNN-GA	0.909	/		
			0.942	/		
P Moeskops et al. [65]	Segment MR brain images into tissue classes	Self-created	/	0.91		
S Erdogmus et al. [80]	Brain extraction	Auto-context CNN	/	0.9773		
K Chen et al. [44]	Segment lesion	Multi-scale CNN	/	0.66		
H Anwar et al. [32]	Segment Glioma tumors in brain	Self-created	/	0.86	0.87	0.90
M and S [62]	Segment deep brain structure	M-net	/	0.8578		

(Dss = Dice Similarity Score; Acc = Accuracy; Com = Complete; Co = Core; E = Enhancing)

\* (TP-CNN is a combination of CSGP+MRILP+MRIGP)

combination of FCN and the new CNN. An excellent result in lung nodule characterization was achieved by Tumor Net in the work of Hussein, Gillies [33]. The network they proposed had nearly 6% improvement in accuracy when compared to methods without utilizing CNN on the same dataset.

### 5.3 Brain

In the application field of brain, CNN has also made outstanding contributions as illustrated in Table 4. It is well known that there are huge differences in appearance between brain tumors, which undoubtedly poses difficult challenges for pathologists or doctors to accurately segment and identify brain tumors. Thus, making use of deep learning to segment and identify brain tumors becomes a promising way of thinking in human medicine [35–42, 97, 99].

Khan, Shah [46] and his team proposed a fully automated CNN-based method for segmenting brain tumors that required a combination of handcrafted features, and achieved excellent results in terms of complete, core, and enhancing tumor. More recently, a fine-tuned VGG-19 (FT-VGG-19) approach to classifying multistage brain tumors that can help radiologists better analyze MRI has yielded compelling results, according to Sajjad, Khan [79] and his team. This classification system not only made use of deep learning to segment the tumor region in the image, but also made usage of data expansion to solve the problem of insufficient data, which showed improved performance in comparison to other existing methods. In addition, to help doctors detect gliomas at an early stage, Anaraki, Ayati [3] and his team came up with a glioma classifier that combines genetic algorithms with a CNN model (CNN-GA) and gets it right more than 90% of the time. This classifier is very flexible because it can be effectively classified in multiple case studies on the classification of brain tumors. Moeskops, Viergever [65] introduced an automatic segmentation method based on CNN, which aims to segment MR brain images into lots of tissue classes. This method tried to learn to recognize the significant information for classification on the basis of training data rather than just dependent on explicit features. The team applied the methods to 5 datasets and



achieved good Dice similarity score (Dss) over segmented tissue classes with the highest result of 0.91, which proves that this automatic segmentation method using CNN has robustness to datasets of different age and acquisition protocol. A brain extraction tool based on the auto-context CNN invented by Salehi, Erdogmus [80] outperforms several results recently reported in brain analysis research. It was superior than other brain extraction methods not only in high Dss results, but also due to its characteristic of registration-free, which means this CNN technique relied less on accuracy of alignment between atlases and anatomy in brain and did not require making assumptions in geometry of brain image. The team of Kamnitsas, Chen [44] introduced a multi-scale CNN solution to segment brain lesion and achieved better performance than previous method on a small-scale dataset. Hussain, Anwar [32] presented several architectures based on CNN to segment Glioma tumors in brain and reaches higher results when compared to state-of-art methods. Mehta and Sivaswamy [62] introduced M-net to help segmenting human brain structure. M-net attracted attention in medical research field as it ran 3 times faster in speed when compared to traditional methods.

Overall, we can see that CNN does bring excellent performance when compared to other state-of-art methods in various areas of medical image analysis. It is often possible to quickly achieve accuracy of up to 90% or more than 95% from data sets that human pathologists may expend considerable effort and time on. There is no doubt that CNN will be of great help to the diagnosis in the medical image analysis area.

## 6 Conclusion

This review presents a comprehensive view on Convolutional Neural Network in the medical field. The content including its developing history, basic operating mechanism, basic structure, classic architecture and applications in medical image analysis.

We first outlined the basic operating mechanism of convolutional neural network, with the survey in neuron, activation function, loss function, gradient descent, feedforward and feedback propagation. And then detailedly introduced 4 categories of convolutional layer, 3 categories of pooling layer, and fully connected layer. The major benefits and limitations of several classic architecture: AlexNet, VGG-Net, Inception, Resnet, U-Net, R-CNN Series, YOLO and SSD are concluded. The applications of CNN among breast, lung and brain area are finally considered from summarizing a number of relevant core journals. We hope that by reading our review, readers can understand the basic principles and popular use framework of CNN, be able to utilize popular format of convolutional or pooling layer like transposed convolutional layer, separable convolutional layer, stochastic pooling in appropriate circumstances, and feel the rapid development of CNN and outstanding contributions brought by CNN in the field of medicine.

A growing number of cases prove that CNN can explore and discover the mysteries of some unknown medical fields. It is expected that the CNN could effectively has the ability to aid diagnostic applications for evaluating more and more diseases in humans in the future.

**Acknowledgements** This work was partially supported by Royal Society International Exchanges Cost Share Award, UK (RP202G0230); Medical Research Council Confidence in Concept Award, UK (MC\_PC\_17171); Hope Foundation for Cancer Research, UK (RM60G0680); British Heart Foundation Accelerator Award, UK; Guangxi Key Laboratory of Trusted Software (kx201901); Fundamental Research Funds for the Central Universities (CDLS-2020-03); Key Laboratory of Child Development and Learning Science (Southeast University), Ministry of Education.

## References

1. Agarap AF (2018) Deep learning using rectified linear units (relu). arXiv preprint arXiv:1803.08375
2. Alakwaa W, Nassef M, Badr A (2017) Lung cancer detection and classification with 3D convolutional neural network (3D-CNN). *Lung Cancer* 8(8):409
3. Anaraki AK, Ayati M, Kazemi F (2019) Magnetic resonance imaging-based brain tumor grades classification and grading via convolutional neural networks and genetic algorithms. *Biocybernetics Biomed Eng* 39(1):63–74
4. Anthimopoulos M, Christodoulidis S, Ebner L, Christe A, Mougiakakou S (2016) Lung pattern classification for interstitial lung diseases using a deep convolutional neural network. *IEEE Trans Med Imaging* 35(5):1207–1216
5. Azamjah N, Soltan-Zadeh Y, Zayeri F (2019) Global trend of breast Cancer mortality rate: a 25-year study. *Asian Pacific J Cancer Prevent: APJCP* 20(7):2015–2020
6. Baldassi C, Malatesta EM, Zecchina R (2019) Properties of the geometry of solutions and capacity of multilayer neural networks with rectified linear unit activations. *Phys Rev Lett* 123(17):6
7. Bardou D, Zhang K, Ahmad SM (2018) Classification of breast cancer based on histology images using convolutional neural networks. *IEEE Access* 6:24680–24693
8. Baudat G, Anouar F (2000) Generalized discriminant analysis using a kernel approach. *Neural Comput* 12(10):2385–2404
9. Bengio Y, Simard P, Frasconi P (1994) Learning long-term dependencies with gradient descent is difficult. *IEEE Trans Neural Netw* 5(2):157–166
10. Bottou L (2010) Large-scale machine learning with stochastic gradient descent, in *Proceedings of COMPSTAT'2010*. Springer, p 177–186
11. Bottou L (2012) Stochastic gradient descent tricks, in *Neural networks: Tricks of the trade*. Springer, p 421–436
12. Bray F et al Erratum: Global cancer statistics (2018): GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers in 185 countries
13. Che Z-G, Chiang T-A, Che Z-H (2011) Feed-forward neural networks training: a comparison between genetic algorithm and back-propagation learning algorithm. *Int J Innovativ Comput Inf Control* 7(10): 5839–5850
14. Chen L-C et al (2018) Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*
15. Chollet F (2017) Xception: Deep learning with depthwise separable convolutions. in *Proceedings of the IEEE conference on computer vision and pattern recognition*
16. Chung H, Lee SJ, Park JG (2016) Deep neural network using trainable activation functions. In *2016 International Joint Conference on Neural Networks (IJCNN)*. IEEE
17. Dasgupta A, Singh S (2017) A fully convolutional neural network based structured prediction approach towards the retinal vessel segmentation. in *2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017)*. IEEE
18. de Koning HJ, van der Aalst CM, de Jong PA, Scholten ET, Nackaerts K, Heuvelmans MA, Lammers JWJ, Weenink C, Yousaf-Khan U, Horeweg N, van 't Westeinde S, Prokop M, Mali WP, Mohamed Hoesein FAA, van Ooijen PMA, Aerts JGJV, den Bakker MA, Thunnissen E, Verschakelen J, Vliegenthart R, Walter JE, ten Haaf K, Groen HJM, Oudkerk M (2020) Reduced lung-cancer mortality with volume CT screening in a randomized trial. *N Engl J Med* 382(6):503–513
19. Du J (2018) Understanding of object detection based on CNN family and YOLO. In *Journal of Physics: Conference Series*. IOP Publishing
20. Dumoulin V, Visin F (2016) A guide to convolution arithmetic for deep learning. arXiv preprint arXiv: 1603.07285
21. Gao F, Wu T, Li J, Zheng B, Ruan L, Shang D, Patel B (2018) SD-CNN: a shallow-deep CNN for improved breast cancer diagnosis. *Comput Med Imaging Graph* 70:53–62
22. Gauen K et al (2017) Low-power image recognition challenge. In *2017 22nd Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE.
23. Girshick R (2015) Fast r-cnn. in *Proceedings of the IEEE international conference on computer vision*
24. Gkioxari G, Girshick R, Malik J (2015) Contextual action recognition with r\* cnn. in *Proceedings of the IEEE international conference on computer vision*
25. Goodfellow I, Bengio Y, Courville A (2016) *Deep Learning*. [SI]. MIT press: 1 Rogers street in Cambridge, Ma 02142

26. Hamidian S et al (2017) 3D convolutional neural network for automatic detection of lung nodules in chest CT. In: Medical Imaging 2017: Computer-Aided Diagnosis. Int Soc Optics Photonics
27. He K et al (2016) Identity mappings in deep residual networks. In European conference on computer vision. Springer
28. Hinton G, Srivastava N, Swersky K (2012) Overview of mini-batch gradient descent. Neural Netw Machine Learning 575
29. Hochreiter S, Schmidhuber J (1997) Long short-term memory. Neural Comput 9(8):1735–1780
30. Howard AG (2013) Some improvements on deep convolutional neural network based image classification. arXiv preprint arXiv:1312.5402
31. Huang J-T, Li J, Gong Y (2015) An analysis of convolutional neural networks for speech recognition. In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE
32. Hussain S, Anwar SM, Majid M (2018) Segmentation of glioma tumors in brain using deep convolutional neural network. Neurocomputing 282:248–261
33. Hussein S et al (2017) Tumomet: Lung nodule characterization using multi-view convolutional neural network with gaussian process. In 2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017). IEEE
34. Iandola FN et al (2016) SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size. arXiv preprint arXiv:1602.07360
35. Jiang YZ et al (2015) Collaborative fuzzy clustering from multiple weighted views. IEEE Trans Cybernetics 45(4):688–701
36. Jiang YZ et al (2017) Seizure classification from EEG signals using transfer Learning, semi-supervised Learning and TSK fuzzy system. IEEE Trans Neural Syst Rehab Eng 25(12):2270–2284
37. Jiang YZ et al (2017) Recognition of epileptic EEG signals using a novel multiview TSK fuzzy system. IEEE Trans Fuzzy Syst 25(1):3–20
38. Jiang Y, Zhao K, Xia K, Xue J, Zhou L, Ding Y, Qian P (2019) A novel distributed multitask fuzzy clustering algorithm for automatic MR brain image segmentation. J Med Syst 43(5):118
39. Jiang YZ et al (2019) Seizure recognition using a novel multitask radial basis function neural network. J Med Imaging Health Informatics 9(9):1865–1870
40. Jiang YZ et al (2019) A novel synthetic CT generation method using multitask maximum entropy clustering. Ieee Access 7:119644–119653
41. Jiang Y, Gu X, Wu D, Hang W, Xue J, Qiu S, Chin-Teng L (2020) A novel negative-transfer-resistant fuzzy clustering model with a shared cross-domain transfer latent space and its application to brain CT image segmentation. IEEE/ACM Trans Comput Biol Bioinform 1. doi: <https://doi.org/10.1109/TCBB.2019.2963873>
42. Jiang Y, Zhang Y, Lin C, Wu D, Lin CT (2020) EEG-based driver drowsiness estimation using an online multi-view and transfer TSK fuzzy system. IEEE Trans Intell Transp Syst: 1–13 <https://doi.org/10.1109/TITS.2020.2973673>
43. Juefei-Xu F, Naresh Boddeti V, Savvides M (2017) Local binary convolutional neural networks. in Proceedings of the IEEE conference on computer vision and pattern recognition
44. Kamnitsas K et al (2015) Multi-scale 3D convolutional neural networks for lesion segmentation in brain MRI. Ischemic Stroke Lesion Segment 13:46
45. Kang G et al (2017) 3D multi-view convolutional neural networks for lung nodule classification. PloS one 12(11):1–21
46. Khan H, Shah PM, Shah MA, Islam S, Rodrigues JJPC (2020) Cascading handcrafted features and convolutional neural network for IoT-enabled brain tumor segmentation. Comput Commun 153:196–207
47. Khirirat S, Feyzmahdavian HR, Johansson M (2017) Mini-batch gradient descent: Faster convergence under data sparsity. in 2017 IEEE 56th Annual Conference on Decision and Control (CDC). IEEE
48. Krizhevsky A, Sutskever I, Hinton GE (2012) ImageNet classification with deep convolutional neural networks. in Advances in neural information processing systems
49. Lawrence S, Giles CL, Ah Chung Tsoi, Back AD (1997) Face recognition: a convolutional neural-network approach. IEEE Trans Neural Netw 8(1):98–113
50. LeCun Y, Boser B, Denker JS, Henderson D, Howard RE, Hubbard W, Jackel LD (1989) Backpropagation applied to handwritten zip code recognition. Neural Comput 1(4):541–551
51. LeCun Y, Bengio Y, Hinton G (2015) *Deep learning*. Nature 521(7553):436–444
52. Lemaréchal C (2012) Cauchy and the gradient method. Doc Math Extra 251:254
53. Lenc K, Vedaldi A (2015) R-cnn minus r. arXiv preprint arXiv:1506.06981
54. Lévy D, Jain A (2016) Breast mass classification from mammograms using deep convolutional neural networks. arXiv preprint arXiv:1612.00542
55. Li J et al (2017) A comparison of deep learning methods for environmental sound detection. In 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE

56. Li Y, Zhang X, Chen D (2018) Csmnet: Dilated convolutional neural networks for understanding the highly congested scenes. In: Proceedings of the IEEE conference on computer vision and pattern recognition
57. Lin M, Chen Q, Yan S (2013) Network in network. arXiv preprint arXiv:1312.4400
58. Liu W et al (2016) Ssd: Single shot multibox detector. In: European conference on computer vision. Springer
59. Marsalli M (2006) Mcculloch-pitts neurons. in The 2008 Annual Meeting of the consortium on cognitive science instruction (ccsi)
60. Mateen M et al (2019) Fundus image classification using VGG-19 architecture with PCA and SVD. *Symmetry* 11(1):1
61. McCulloch WS, Pitts W (1943) A logical calculus of the ideas immanent in nervous activity. *Bull Mathematical Biophys* 5(4):115–133
62. Mehta R, Sivaswamy J (2017) M-net: A convolutional neural network for deep brain structure segmentation. In 2017 IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017). IEEE.
63. Mikolov T, et al (2010) Recurrent neural network based language model. in Eleventh annual conference of the international speech communication association
64. Minsky M, Papert SA (2017) Perceptrons: An introduction to computational geometry. MIT press
65. Moeskops P, Viergever MA, Mendrik AM, de Vries LS, Benders MJNL, Išgum I (2016) Automatic segmentation of MR brain images with a convolutional neural network. *IEEE Trans Med Imaging* 35(5): 1252–1261
66. Mordang J-J et al (2016) Automatic microcalcification detection in multi-vendor mammography using convolutional neural networks. In: International Workshop on Breast Imaging. Springer
67. Nielsen MA (2015) Neural networks and deep learning, vol 2018. Determination press San Francisco, CA, USA
68. Pascanu R, Mikolov T, Bengio Y (2012) Understanding the exploding gradient problem. *CoRR*, abs/1211.5063. 2:417
69. Pedamonti D (2018) Comparison of non-linear activation functions for deep neural networks on MNIST classification task. arXiv preprint arXiv:1804.02763
70. Purkait P, Zhao C, Zach C (2017) SPP-Net: Deep absolute pose regression with synthetic views. arXiv preprint arXiv:1712.03452
71. Rastegari M, et al (2016) Xnor-net: Imagenet classification using binary convolutional neural networks. In European conference on computer vision. Springer
72. Reddy S (2018) Use of artificial intelligence in healthcare delivery, in eHealth-Making Health Care Smarter. *IntechOpen*:81–97
73. Ren S et al (2015) Faster r-cnn: Towards real-time object detection with region proposal networks. in Advances in neural information processing systems
74. Ronneberger O, Fischer P, Brox T (2015) U-net: Convolutional networks for biomedical image segmentation. In International Conference on Medical image computing and computer-assisted intervention. Springer
75. Rosenblatt F (1958) The perceptron: a probabilistic model for information storage and organization in the brain. *Psychol Rev* 65(6):386–408
76. Rouhi R, Jafari M, Kasaei S, Keshavarzian P (2015) Benign and malignant breast tumors classification based on region growing and CNN segmentation. *Expert Syst Appl* 42(3):990–1002
77. Ruder S (2016) An overview of gradient descent optimization algorithms. arXiv preprint arXiv:1609.04747
78. Rumelhart DE, Hinton GE, Williams RJ (1986) Learning representations by back-propagating errors. *Nature* 323(6088):533–536
79. Sajjad M, Khan S, Muhammad K, Wu W, Ullah A, Baik SW (2019) Multi-grade brain tumor classification using deep CNN with extensive data augmentation. *J Comput Sci* 30:174–182
80. Salehi SSM, Erdogmus D, Gholipour A (2017) Auto-context convolutional neural network (auto-net) for brain extraction in magnetic resonance imaging. *IEEE Trans Med Imaging* 36(11):2319–2330
81. Salekin MS, Jelodar AB, Kushol R (2019) Cooking state recognition from images using inception architecture. In 2019 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST). IEEE
82. Scardapane S, van Vaerenbergh S, Totaro S, Uncini A (2019) Kafnets: kernel-based non-parametric activation functions for neural networks. *Neural Netw* 110:19–32
83. Schairer C, Mink PJ, Carroll L, Devesa SS (2004) Probabilities of death from breast cancer and other causes among female breast cancer patients. *J Natl Cancer Inst* 96(17):1311–1321
84. Schmidhuber J (2015) Deep learning in neural networks: an overview. *Neural Netw* 61:85–117
85. Schmidt WF et al (1993) Initializations, back-propagation and generalization of feed-forward classifiers. In IEEE International Conference on Neural Networks. IEEE
86. Schuster M, Paliwal KK (1997) Bidirectional recurrent neural networks. *IEEE Trans Signal Process* 45(11):2673–2681

87. Shen W et al (2015) Multi-scale convolutional neural networks for lung nodule classification. In: International Conference on Information Processing in Medical Imaging. Springer
88. Sheng T et al (2018) A quantization-friendly separable convolution for mobilenets. In 2018 1st Workshop on Energy Efficient Machine Learning and Cognitive Computing for Embedded Applications (EMC2). IEEE
89. Simonyan K, Zisserman A (2014) Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556
90. Su H et al (2015) Region segmentation in histopathological breast cancer images using deep convolutional neural network. In: 2015 IEEE 12th International Symposium on Biomedical Imaging (ISBI). IEEE
91. Sun M, Song Z, Jiang X, Pan J, Pang Y (2017) Learning pooling for convolutional neural network. *Neurocomputing* 224:96–104
92. Sun W, Tseng TL(B), Zhang J, Qian W (2017) Enhancing deep convolutional neural network scheme for breast cancer diagnosis with unlabeled data. *Comput Med Imaging Graph* 57:4–9
93. Szegedy C et al (2015) Going deeper with convolutions. In Proceedings of the IEEE conference on computer vision and pattern recognition
94. Szegedy C et al (2016) Rethinking the inception architecture for computer vision. in Proceedings of the IEEE conference on computer vision and pattern recognition
95. Szegedy C et al (2017) Inception-v4, inception-resnet and the impact of residual connections on learning. in Thirty-first AAAI conference on artificial intelligence
96. Ting FF, Tan YJ, Sim KS (2019) Convolutional neural network improvement for breast cancer classification. *Expert Syst Appl* 120:103–115
97. Toğaçar M, Ergen B, Cömert Z (2020) BrainMRNet: brain tumor detection using magnetic resonance images with a novel convolutional neural network model. *Med Hypotheses* 134:109531
98. Waks AG, Winer EP (2019) Breast cancer treatment: a review. *Jama* 321(3):288–300
99. Xue Y, Farhat FG, Boukrina O, Barrett AM, Binder JR, Roshan UW, Graves WW (2020) A multi-path 2.5 dimensional convolutional neural network system for segmenting stroke lesions in brain MRI images. *NeuroImage: Clin* 25:102118
100. Yu F, Koltun V (2015) Multi-scale context aggregation by dilated convolutions. arXiv preprint arXiv:1511.07122
101. Yu D et al (2014) Mixed pooling for convolutional neural networks. In International conference on rough sets and knowledge technology. Springer
102. Yuan Z-W, Zhang J (2016) Feature extraction and image retrieval based on AlexNet. in Eighth International Conference on Digital Image Processing (ICDIP 2016). International Society for Optics and Photonics
103. Zarándy Á et al (2015) Overview of CNN research: 25 years history and the current trends. In 2015 IEEE International Symposium on Circuits and Systems (ISCAS). IEEE
104. Zeiler MD, Fergus R (2013) Stochastic pooling for regularization of deep convolutional neural networks. arXiv preprint arXiv:1301.3557
105. Zeiler MD, Fergus R (2014) Visualizing and understanding convolutional networks. In European conference on computer vision. Springer
106. Zeiler MD, Taylor GW, Fergus R (2011) Adaptive deconvolutional networks for mid and high level feature learning. In 2011 International Conference on Computer Vision. IEEE
107. Zhang Z, Sarhadi M (1993) A modified neuron activation function which enables single layer perceptrons to solve some linearly inseparable problems. In Proceedings of 1993 International Conference on Neural Networks (IJCNN-93-Nagoya, Japan). IEEE
108. Zubair S, Yan F, Wang W (2013) Dictionary learning based sparse coefficients for audio classification with max and average pooling. *Digital Signal Process* 23(3):960–970

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## Affiliations

Xujing Yao<sup>1</sup> · Xinyue Wang<sup>1</sup> · Shui-Hua Wang<sup>1,2,3</sup> · Yu-Dong Zhang<sup>1,2</sup>

<sup>1</sup> School of Informatics, University of Leicester, Leicester LE1 7RH, UK

<sup>2</sup> Department of Information Systems, Faculty of Computing and Information Technology, King Abdulaziz University, Jeddah 21589, Saudi Arabia

<sup>3</sup> School of Architecture Building and Civil Engineering, Loughborough University, Loughborough LE11 3TU, UK