

# Machine learning

BY YOUJUN HU

yjhu@ipp.cas.cn

## 1 Introduction

Artificial intelligence (AI) research has tried many different approaches since its founding. In the first decades of the 21st century, the AI research is dominated by highly mathematical optimization machine learning (ML), which has proved successful in solving many challenging real life problems.

Many problems in AI can be solved theoretically by searching through many possible solutions: Reasoning can be reduced to performing a search. Simple exhaustive searches are rarely sufficient for most real-world problems. The solution, for many problems, is to use "heuristics" or "rules of thumb" that prioritize choices in favor of those more likely to reach a goal. A very different kind of search came to prominence in the 1990s, based on the mathematical theory of optimization. Modern machine learning is based on these methods. Instead, of using detailed explanations to guide the search, it uses a combination of [1]: (a) general architectures; (b) trying trillions of possibilities, guided by simple ideas (like gradient descent) for improvement; and (c) the ability to recognize progress (by defining a objective function).

I am interested in applying machine learning to problems in computational physics problems that traditional numerical methods can not easily handle either because of its computational costs being too high or its traditional algorithms are too complicated to easily implement.

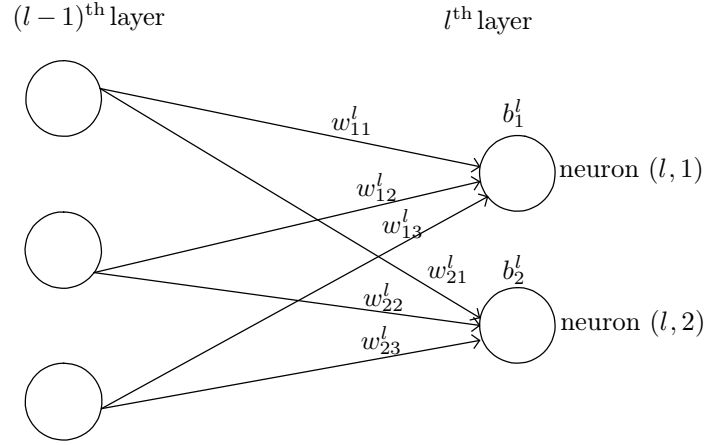
Enrico Fermi once criticized the complexity of a model (that contains many free parameters) by quoting Johnny von Neumann "With four parameters I can fit an elephant, and with five I can make him wiggle his trunk". What Fermi implies is that it is easy to fit existing data and what is important is to have a model with predicting capability (fitting data not seen yet). The artificial neural network method tackles this difficulty by increasing the number of free parameters to millions, with the hope of obtaining predicting capability.

## 2 Neural network

Neural networks consists of multiple layers of interconnected nodes (neurons), each having a weight for a connection, a bias, and an activation function. Each layer build upon the previous layer. This progression of computations through the network is called forward propagation. Another process called backpropagation uses algorithms which moves backwards through the layers to efficiently compute the partial derivatives of the objective function with respect to the weights and biases. Combining the forward and backward propagation, we can calculate errors in predictions and then adjusts the weights and biases using the gradient descent method. This process is called training/learning.

### 2.1 Node (neuron or unit), weight, bias, and activation

As is shown in Fig. 1, we use  $w_{jk}^l$  to denote the weight for the connection from the  $k^{\text{th}}$  neuron in the  $(l-1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer. Use  $b_j^l$  to denote the bias of the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer.



**Figure 1.** Definition of layers, neurons, weights, and biases in a neural network. The  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer is referred to as neuron  $(l, j)$

We use  $a_j^l$  to denote the output (activation) of the  $j^{\text{th}}$  neuron in  $l^{\text{th}}$  layer. A neural network model assumes that  $a_j^l$  is related to the  $a^{l-1}$  (output of the previous layer) via

$$a_j^l = \sigma \left( \sum_k w_{jk}^l a_k^{l-1} + b_j^l \right), \quad (1)$$

where the summation is over all neurons in the  $(l-1)^{\text{th}}$  layer and  $\sigma$  is a function called activation function which can take various forms, e.g., step function,

$$\sigma(z) = \begin{cases} 1 & \text{if } z \geq 0 \\ 0 & \text{else} \end{cases}, \quad (2)$$

rectified linear unit (ReLU),

$$\sigma(z) = \max(0, z), \quad (3)$$

and sigmoid (“S”-shaped curve, also called logistic function)

$$\sigma(z) = \frac{1}{1 + \exp(-z)}. \quad (4)$$

For notation ease, define  $z_j^l$  by

$$z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l, \quad (5)$$

which can be interpreted as an weighted input to the neuron  $(l, j)$ , then Eq. (1) is written as

$$a_j^l = \sigma(z_j^l). \quad (6)$$

In matrix form, Eq. (5) is written as

$$z^l = w^l a^{l-1} + b^l, \quad (7)$$

where  $w^l$  is a  $J \times K$  matrix,  $z^l$  and  $b^l$  are column vectors of length  $J$ ,  $a^{l-1}$  is a column vector of length  $K$ , where  $J$  and  $K$  are the number of neurons in the  $l^{\text{th}}$  layer and  $(l-1)^{\text{th}}$  layer, respectively.

The input layer is where data inputs are provided, and the output layer is where the final prediction is made. The input and output layers of a deep neural network are called visible layers. The layers between the input layer and output layer are called hidden layers. Note that the input layer is usually not considered as a layer of the network since it does not involve any computation. In tensorflow, layers refer to the computing layers (i.e., hidden layers and the output layer, not including the input layer). The activation function of each layer can be different. The activation function of the output layer is often chosen as None, ReLU, logistic/tanh, and is usually different from those used in the hidden layers. Here “None” means activation  $\sigma(z) = z$ .

## 2.2 Objective function

Define an objective function (can be called loss or cost function depending on contexts) by

$$C(w, b) \equiv \frac{1}{2n} \sum_x \|y(x) - a^L\|^2, \quad (8)$$

where  $w$  and  $b$  denotes the collection of all weights and biases in the network,  $n$  is the total number of training examples  $x$ , the summation is over all the training examples,  $y(x)$  is the desired output from the network (i.e., correct answer) when  $x$  is the input, and  $a^L$  is the actual output from the output layer of the network and is a function of  $w, b$ , and  $x$ . Note that  $y$  and  $a^L$  are vectors (with number of elements being the number of neurons in the output layer) and  $\|\dots\|$  denotes the vector norm. Explicitly writing out the vector norm, Eq. (8) is written as

$$C(w, b) \equiv \frac{1}{2n} \sum_x \sum_{j=1}^{N_L} (y_j(x) - a_j^L)^2, \quad (9)$$

where  $N_L$  is the number of neurons in the output layer.

The cost function is the average error of the approximate solution away from the desired exact solution. The goal of a learning algorithm is to find weights and biases that minimize the cost function. A method of minimizing the cost function over  $(w, b)$  is the gradient descent method:

$$w_{jk}^l \rightarrow w_{jk}^l - \eta \frac{\partial C}{\partial w_{jk}^l}, \quad (10)$$

$$b_j^l \rightarrow b_j^l - \eta \frac{\partial C}{\partial b_j^l}, \quad (11)$$

where  $\eta$  is called learning rate, which should be positive.

In using the gradient descent method, we need to compute the partial derivatives  $\partial C / \partial w_{jk}^l$  and  $\partial C / \partial b_j^l$ . Next we will discuss how to compute them.

## 2.3 Gradients of objective function

Note that the loss function involves an average over all the training examples. Denote the loss function for a specific training example by  $C_x$ , i.e.,

$$C_x = \frac{1}{2} \sum_{j=1}^{N_L} (y_j(x) - a_j^L)^2, \quad (12)$$

then expression (9) is written as

$$C = \frac{1}{n} \sum_x C_x, \quad (13)$$

Then the partial derivatives  $\partial C / \partial w_{jk}^l$  and  $\partial C / \partial b_j^l$  can be written as the sum of  $\partial C_x / \partial w_{jk}^l$  and  $\partial C_x / \partial b_j^l$ , i.e.,

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{1}{n} \sum_x \frac{\partial C_x}{\partial w_{jk}^l}, \quad (14)$$

$$\frac{\partial C}{\partial b_j^l} = \frac{1}{n} \sum_x \frac{\partial C_x}{\partial b_j^l}. \quad (15)$$

The above formulas indicate that, once  $\partial C_x / \partial w_{jk}^l$  and  $\partial C_x / \partial b_j^l$  are known, obtaining  $\partial C / \partial w_{jk}^l$  and  $\partial C / \partial b_j^l$  is trivial, i.e., just averaging them. Therefore, we will focus on  $C_x$  (i.e., the cost function for a fixed training example) and discuss how to compute  $\partial C_x / \partial w_{jk}^l$  and  $\partial C_x / \partial b_j^l$ .

In practice, we do not sum over all the training examples. Instead, we average the derivative over a small number (say 16) of training examples (a mini batch) and use these approximate derivatives to advance a step. For the next step, we stochastically change to using a different mini batch. This is called stochastic gradient descent (SGD) method.

## 2.4 Back-propagating algorithm

The cost function  $C_x$  is a function of weights and biases of all neurons (the input  $x$  and output  $y(x)$  are fixed parameters). For a specific neuron  $(l, j)$ , its weights and biases enter  $C_x$  via the combination  $z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$ . Then it is useful to define the following partial derivative:

$$\delta_j^l \equiv \frac{\partial C_x}{\partial z_j^l}, \quad (16)$$

where the partial derivative are taken with fixed weights and biases for all neurons except neuron  $(l, j)$ . Note that the  $a_k^{l-1}$  appearing in the expression of  $z_j^l$  does not depend on  $w_{jk}^l$  or  $b_j^l$ . It only depends on the weights and biases in the layers  $\leq (l-1)$ , which are all fixed when taking the derivative in expression (16).  $\delta_j^l$  defined in expression (16) is often called the error of neuron  $(l, j)$ .

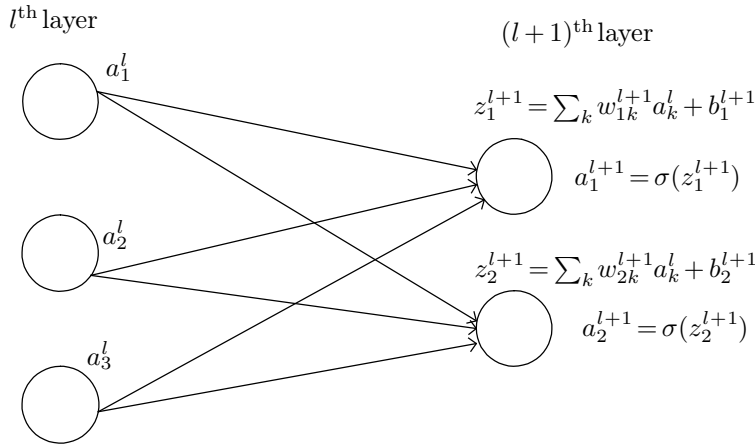
Using the chain rule,  $\partial C_x / \partial w_{jk}^l$  and  $\partial C_x / \partial b_j^l$  can be expressed in terms of  $\delta_j^l$ :

$$\frac{\partial C_x}{\partial b_j^l} = \frac{\partial C_x}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} = \delta_j^l, \quad (17)$$

and

$$\frac{\partial C_x}{\partial w_{jk}^l} = \frac{\partial C_x}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} = \delta_j^l a_k^{l-1}. \quad (18)$$

Therefore, if  $\delta_j^l$  is known, it is trivial to compute the gradients needed in the gradient descent method.



For the output layer ( $L^{\text{th}}$  layer),  $\delta_j^L$  defined in Eq. (16) is written as

$$\delta_j^L = \frac{\partial C_x}{\partial z_j^L} = \frac{\partial C_x}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L}. \quad (19)$$

The dependence of  $C_x$  on  $a_j^L$  is explicitly given by Eq. (12), from which the above expression for  $\delta_j^L$  is written as

$$\delta_j^L = (a_j^L - y_j(x)) \sigma'(z_j^L). \quad (20)$$

Therefore  $\delta_j^L$  is easy to compute.

Backpropagation is a way of computing  $\delta_j^l$  for every layer using recurrence relations: the relation between  $\delta^l$  and  $\delta^{l+1}$ . Noting how the error is propagating through the network, we know the following identity:

$$\frac{\partial C_x}{\partial z_j^l} dz_j^l = \sum_j \frac{\partial C_x}{\partial z_j^{l+1}} dz_j^{l+1}, \quad (21)$$

with

$$dz_j^{l+1} = w_{jJ}^{l+1} d(a_J^l), \quad (22)$$

i.e.,

$$dz_j^{l+1} = w_{jJ}^{l+1} \sigma'(z_j^l) dz_j^l. \quad (23)$$

Therefore

$$\frac{\partial C_x}{\partial z_j^l} = \sum_j \frac{\partial C_x}{\partial z_j^{l+1}} w_{jJ}^{l+1} \sigma'(z_j^l). \quad (24)$$

i.e.,

$$\delta_j^l = \sum_j \delta_j^{l+1} w_{jJ}^{l+1} \sigma'(z_j^l). \quad (25)$$

Equation (25) gives the recurrence relations of computing  $\delta^l$  from  $\delta^{l+1}$ . This is called the back-propagation algorithm. Eq. (25) can be written in the matrix form:

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l), \quad (26)$$

where  $T$  stands for matrix transpose,  $\odot$  is the element-wise product.

### 3 Automatic differentiation

Automatic differentiation (autodiff) is a set of techniques for computing derivatives of numeric functions expressed as source code (i.e. the internal mechanism of the function is known). It works by breaking down functions into elementary operations (addition, multiplication, etc.) whose derivatives are known, and then applies the chain rule to compute the derivatives.

Autodiff can be considered as a kind of symbolic differentiation. The difference of autodiff from the traditional symbolic differentiation is that the goal is not to get a compact formula for humans to understand, but for computers to evaluate. Therefore the final result from a autodiff is not an analytical formula, but numerical data. This goal also makes autodiff more efficient since autodiff does not need to perform some intermedia processes that appear when you use traditional symbolic differentiation to get a formula and then numerically evaluate the formula.

Autodiff can be better than numerical differentiation (e.g., finite difference) in that it avoid truncation errors. When you are given a black-box function, you can not use autodiff since the internal mechanism of the function is unknown. In this case, the only choice is to use numerical differentiation.

Autodiff has two main modes: forward and backward. Forward mode computes a single derivative during one pass of the expression tree. The back mode computes all the derivatives in a single pass of the expression tree, avoiding some computational repetition (compared using forward mode to compute all the derivatives sperately). It's more efficient for functions with many inputs and few outputs, making it ideal for machine learning where we often compute gradients of scalar loss functions with respect to many parameters. Backpropagation is a specific instance of backward-mode autodiff. Here is a simple Python code implementing autodiff:

Forward:

```
class Expression:
    def __add__(exp1, exp2):
        return Plus(exp1, exp2)
    def __mul__(exp1, exp2):
        return Multiply(exp1,exp2)
class Variable(Expression):
    def __init__(self,value):
        self.value = value
    def evaluate(self):
        return self.value
    def derive(self, v):
        return 1 if self == v else 0

class Plus(Expression):
```

```

def __init__(self,exp1,exp2):
    self.a = exp1
    self.b = exp2
def evaluate(self):
    return self.a.evaluate() + self.b.evaluate()
def derive(self,v):
    return self.a.derive(v) + self.b.derive(v)

class Multiply(Expression):
    def __init__(self,exp1,exp2):
        self.a = exp1
        self.b = exp2
    def evaluate(self):
        return self.a.evaluate() * self.b.evaluate()
    def derive(self, v):
        return (self.a.derive(v)*self.b.evaluate()
                +self.b.derive(v)*self.a.evaluate())

# Example: derivatives of z(x,y) at (x, y) = (2, 3)
x = Variable(2)
y = Variable(3)
z = x * (x + y) + y * y
print(z.derive(x)) # dz/dx, Output: 7
print(z.derive(y)) # dz/dy, Output: 8

```

This simple example illustrates several important concepts:

- \* Class, subclass, inheritance
- \* Operator overloading
- \* Polymorphism
- \* Recursion

The following is the backward (or reverse) method.

```

class Expression:
    def __add__(exp1, exp2):
        return Plus(exp1,exp2)
    def __mul__(exp1, exp2):
        return Multiply(exp1, exp2)

class Variable(Expression):
    def __init__(self,value):
        self.value = value
        self.partial = 0
    def evaluate(self):
        return self.value
    def derive(self,seed):
        self.partial += seed

class Plus(Expression):
    def __init__(self,exp1,exp2):
        self.a = exp1
        self.b = exp2
    def evaluate(self):
        return self.a.evaluate() + self.b.evaluate()
    def derive(self,seed):
        self.a.derive(seed)
        self.b.derive(seed)

```

```

class Multiply(Expression):
    def __init__(self, exp1, exp2):
        self.a = exp1
        self.b = exp2
    def evaluate(self):
        return self.a.evaluate() * self.b.evaluate()
    def derive(self, seed):
        self.a.derive(seed * self.b.evaluate())
        self.b.derive(seed * self.a.evaluate())

```

```

x = Variable(2)
y = Variable(3)
z = x * (x + y) + y * y
z.derive(1)
print(x.partial) # dz/dx Output: 7
print(y.partial) # dz/dy Output: 8

```

Another method is to use dual numbers. Dual numbers are expressions of the form  $a + b\varepsilon$ , where  $a$  and  $b$  are real numbers, and  $\varepsilon$  is a symbol taken to satisfy  $\varepsilon^2 = 0$  with  $\varepsilon \neq 0$ . Using a dual number in the Taylor series, we obtain

$$f(a + b\varepsilon) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (b\varepsilon)^n = f(a) + bf'(a)\varepsilon, \quad (27)$$

since all terms involving  $\varepsilon^2$  and greater powers are zero by the definition of  $\varepsilon$ . We find the coefficient of  $\varepsilon$  in the result is the first derivative  $f'(a)$ . This result can be used in computer programs to find derivative of a function by defining a class and overloading the basic operators, e.g.

$$a_1 + b_1\varepsilon + a_2 + b_2\varepsilon = (a_1 + a_2) + (b_1 + b_2)\varepsilon, \quad (28)$$

$$(a_1 + b_1\varepsilon)(a_2 + b_2\varepsilon) = a_1a_2 + (a_1b_2 + a_2b_1)\varepsilon. \quad (29)$$

Here is a Python code:

```

class Dual:
    def __init__(self, realPart, infs=0):
        self.realPart = realPart
        self.infs = infs

    def __add__(self, other):
        return Dual(
            self.realPart + other.realPart,
            self.infs + other.infs
        )

    def __mul__(self, other):
        return Dual(
            self.realPart * other.realPart,
            other.realPart * self.infs + self.realPart * other.infs
        )

def f(x, y):
    return x * (x + y) + y * y
x = Dual(2)
y = Dual(3)
epsilon = Dual(0, 1)
a = f(x + epsilon, y)
b = f(x, y + epsilon)

```





```

regressor = Linear_Regression()
Nd=11
X = np.array([i for i in range(Nd)])
Y = np.array([2*i for i in range(Nd)]) + np.random.uniform(high=5.0,size=Nd)
fig, ax = plt.subplots()
ax.plot(X,Y, 'k.',label='data')
Y_pred = regressor.predict(X)
ax.plot(X, Y_pred, label='Initial fit line')

learning_rate = 0.01
i = 0
while i<100:
    regressor.update_coeffs(X,Y,learning_rate)
    i = i+1

Y_pred = regressor.predict(X)
ax.plot(X, Y_pred, 'b-',label='Final Fit Line')
ax.legend()
plt.show()

```

The loss function is defined by the mean square error, which is not directly used in the above code. Only the partial derivatives of the loss function is directly used.

## 6 Logistic regression for binary classification

Hypothesis function (the model): Denote the output of the model by  $\hat{y}$ , which is given by

$$\hat{y} = \sigma(z), \quad (34)$$

where  $z = w \cdot x + b$  and  $\sigma$  is the sigmoid function given by Eq. (4). The model is nonlinear in the unknowns  $w$  and  $b$ .

The loss function is chosen as

$$L = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)], \quad (35)$$

where  $y_i$  is the correct answer of the  $i$ th training example ( $y_i$  can take only two values, 0 or 1). The value of  $\hat{y}_i$  is interpreted as the probability of  $y$  being 1.

Because the model function is nonlinear and the loss function is complicated, there is usually no closed-form solution that minimizes the loss function. Iterative methods, such as gradient descent, are needed to solve for  $w$  and  $b$ . The partial derivatives needed in the gradient descent method can be written as

$$\begin{aligned}
\frac{\partial L}{\partial w} &= -\frac{1}{m} \sum_i \left[ y_i \frac{1}{\hat{y}_i} \sigma'(z) x_i - (1 - y_i) \frac{1}{(1 - \hat{y}_i)} \sigma'(z) x_i \right] \\
&= -\frac{1}{m} \sum_i \left\{ \left[ y_i \frac{1}{\hat{y}_i} - (1 - y_i) \frac{1}{(1 - \hat{y}_i)} \right] \sigma'(z) x_i \right\} \\
&= -\frac{1}{m} \sum_i \left\{ \left[ \frac{y_i - \hat{y}_i}{\hat{y}_i(1 - \hat{y}_i)} \right] \sigma'(z) x_i \right\} \\
&= -\frac{1}{m} \sum_i \left\{ \left[ \frac{y_i - \hat{y}_i}{\hat{y}_i(1 - \hat{y}_i)} \right] \sigma(1 - \sigma) x_i \right\} \\
&= -\frac{1}{m} \sum_i \left\{ \left[ \frac{y_i - \hat{y}_i}{\hat{y}_i(1 - \hat{y}_i)} \right] \hat{y}_i(1 - \hat{y}_i) x_i \right\} \\
&= -\frac{1}{m} \sum_i [(y_i - \hat{y}_i) x_i] \quad (36)
\end{aligned}$$

Using

$$\frac{d\sigma}{dz} = \frac{1}{(1 + \exp(-z))^2} \exp(-z) = \sigma^2 \exp(-z) = \sigma(1 - \sigma) \quad (37)$$

The above formula is simplified as

$$\begin{aligned} \frac{\partial L}{\partial w} &= -\frac{1}{m} \sum_i \left\{ \left[ \frac{y_i - \hat{y}_i}{\hat{y}_i(1 - \hat{y}_i)} \right] \sigma(1 - \sigma) x_i \right\} \\ &= -\frac{1}{m} \sum_i \left\{ \left[ \frac{y_i - \hat{y}_i}{\hat{y}_i(1 - \hat{y}_i)} \right] \hat{y}_i(1 - \hat{y}_i) x_i \right\} \\ &= -\frac{1}{m} \sum_i [(y_i - \hat{y}_i) x_i] \end{aligned} \quad (38)$$

Similarly, we obtain

$$\frac{\partial L}{\partial w} = -\frac{1}{m} \sum_i (y_i - \hat{y}_i). \quad (39)$$

## Bibliography

- [1] Michael A Nielsen. *Neural networks and deep learning*, volume 25. Determination Press: <http://neuralnetworksanddeeplearning.com/>, 2015.