

# 基于权值变化的 BP 神经网络自适应学习率改进研究<sup>①</sup>

朱振国, 田松禄

(重庆交通大学 信息科学与工程学院, 重庆 400074)

通讯作者: 田松禄, E-mail: mymailwith163@163.com

**摘要:** 针对传统神经网络的学习率由人为经验性设定, 存在学习率设置过大或过小, 容易导致无法收敛或收敛速度慢的问题, 本文提出基于权值变化的自适应学习率改进方法, 改善传统神经网络学习率受人为经验因素影响的弊端, 提高误差精度, 并结合正态分布模型与梯度上升法, 提高收敛速度. 本文以 BP 神经网络为例, 对比固定学习率的神经网络, 应用经典 XOR 问题仿真验证, 结果表明本文的改进神经网络具有更快的收敛速度和更小的误差.

**关键词:** 神经网络; 自适应学习率; 正态分布模型; 梯度上升法; XOR 问题

引用格式: 朱振国, 田松禄. 基于权值变化的 BP 神经网络自适应学习率改进研究. 计算机系统应用, 2018, 27(7): 205-210. <http://www.c-s-a.org.cn/1003-3254/6410.html>

## Improvement of Learning Rate of Feed Forward Neural Network Based on Weight Gradient

ZHU Zhen-Guo, TIAN Song-Lu

(College of Information Science and Engineering, Chongqing Jiaotong University, Chongqing 400074, China)

**Abstract:** An adaptive learning rate improvement method, based on weight change, is proposed to improve the learning rate of traditional neural network in this study. If the learning rate is too large or too small, neural network is too difficult or too slow to converge. To offset this disadvantage, the study put forward a new learning rate, based on weight gradient, to improve the convergence rate and improve the traditional neural network learning rate affected by the human experienced factors, and combined with normal distribution and gradient rise method, to size up error accuracy and convergence speed. Taking BP neural network as an example, comparing the fixed learning rate neural network, and applying classical XOR problem simulation, we verify the proposed method. The results show that this improved neural network has faster convergence speed and smaller error.

**Key words:** neural network; adaptive learning rate; normal distribution model; method of gradient increase; XOR issue

## 引言

BP 神经网络 (Back-Propagation Neural Network), 是由 Rumelhart 和 McClelland 等科学家提出, 利用输入信号前向传播、误差反馈信号反向传播和梯度下降的原理, 并通过链式求导法则, 获取权值更新变化大小的依据, 使权值可以按照一定的大小进行更新, 达到减小误差、得到理想输出的一种算法. 常用于预测、回归问题的判别, 是目前应用最为广泛的神经网络之一<sup>[1]</sup>.

但传统 BP 神经网络, 比如: 学习率为固定值, 学习率设置偏大, 容易导致学习震荡甚至发散, 而无法收敛; 学习率设置偏小, 容易导致学习速率慢, 收敛过于缓慢; 对于这种由于学习率人为设定不合理的问题, 不能较好地建立输入输出的非线性映射关系, 而导致 BP 神经网络难以推广应用<sup>[2]</sup>.

针对 BP 神经网络学习率的人为设定不合理的问题, 本文提出基于权值变化的自适应学习率模型, 改进

<sup>①</sup> 收稿时间: 2017-10-22; 修改时间: 2017-11-10; 采用时间: 2017-11-27; csa 在线出版时间: 2018-06-27

了传统神经网络的固定学习率设置不合理的弊端;并将正态分布结合神经网络的误差函数,加快收敛速度;利用梯度上升法,以保证正态分布的合理应用.

### 1 神经网络

对于三层 BP 神经网络,输入  $X_1, X_2, X_3 \dots X_n$ , 输出为  $Y$ , 隐含层输入权值为  $W_{ij}^L$ , 输出层权值为  $W_{ij}^{L+1}$ ,  $b_j$  为阈值,  $L$  为层数,  $ij$  表示前层第  $i$  个的和后层第  $j$  个神经元,  $f(\cdot)$  表示激活函数<sup>[3]</sup>. 隐层神经元净输入值为:

$$S = W_{ij}^L \times X_n + b_j^L \tag{1}$$

激活函数  $f(\cdot)$  采用 Sigmoid 函数, 激活后的值域为  $(0, 1)$ , 即:

$$f(S_j^L) = X_{ij}^L = \frac{1}{1 + e^{-S_j^L}} \tag{2}$$

$f(\cdot)$  的导数为:

$$f(S_j^L)' = f(S_j^L) \times (1 - f(S_j^L)) \tag{3}$$

期望输出用  $d$  表示, 实际输出为  $Y$ ; 误差函数  $err$  的表达式为<sup>[4]</sup>:

$$err = \sum \|d - Y\|_2 \tag{4}$$

式 (4) 可以看出, 存在理想极小值点  $err=0$ , 但实际很难达到该点, 通常是根据误差反向传播与梯度下降法<sup>[5]</sup>, 多次迭代更新权值, 使实际输出  $Y$  无限逼近期望输出  $d$ , 达到误差  $err$  逼近 0 的目的<sup>[6-9]</sup>.

## 2 正态分布模型和自适应学习率的 BP 神经网络

### 2.1 正态分布模型

引入正态分布模型到 BP 神经网络中, 将误差  $err$  作为正态分布函数的自变量, 令正态分布模型的期望  $u$  为 0, 正态分布函数值取得最大值, 误差  $err$  趋近于  $u$ , 如图 1.

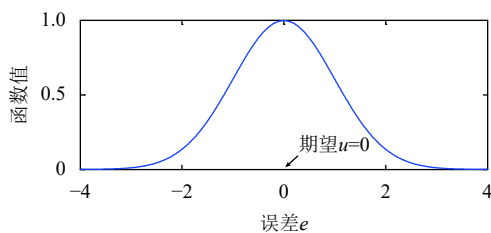


图 1 正态分布函数图

当期望  $u=0$  时, 正态分布函数:

$$g = \frac{1}{\delta \times \sqrt{2\pi}} \times e^{-\left(\frac{err}{\delta}\right)^2} \tag{5}$$

网络训练目标是取得正态分布函数最大值, 目标达成, 则网络误差为 0, 权值更新达到最佳状态.

本文借鉴用于取得局部最小值的梯度下降法 (要求误差函数为凹函数) 思想, 反向推理, 采用梯度上升法 (要求误差函数为凸函数) 寻找正态分布的最大值.

以图 1 和式 (5) 为例, 解释梯度上升法能取得局部最大值的原理:

$$en = err + g' \times \alpha \tag{6}$$

当  $err < 0, g' > 0$  时,  $en > err$ ,

当  $err > 0, g' < 0$  时,  $en < err$ ,

当  $err = 0, g' = 0$  时,  $en = err$ ; 此时  $g$  取得最大值.

其中,  $en$  为  $err$  更新后的误差值,  $\alpha$  为学习率, 值域为  $(0, 1)$ .

### 2.2 自适应学习率模型

提出基于权值变化的自适应学习率定义为:

$$\beta(t) = \frac{2}{1 + e^{-|t| \times 10^n}} - 1 \tag{7}$$

其中  $t$  是 BP 神经网络的权值变化:

$$t = \frac{\partial err}{\partial W_j^L} = \frac{\partial err}{\partial f(S_j^L)} \times \frac{\partial f(S_j^L)}{\partial S_j^L} \times \frac{\partial S_j^L}{\partial W_j^L} \tag{8}$$

$\beta(t)$  函数曲线图为图 2.

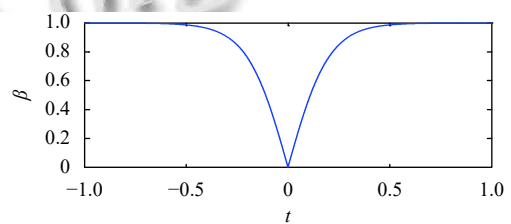


图 2 自适应学习率的曲线图

图 2 中  $n$  为  $\beta$  的倾斜参数. 由式 (7)、式 (8) 及图 2 可以得出结论, 当训练接近理想时, 权值的变化  $t$  趋于极小的值, 此时学习率  $\beta$  也是一个极小的值, 不利于训练的进行, 于是用  $10^n$  扩大  $t$  的值, 调整  $\beta$  函数对  $t$  的敏感程度.

BP 神经网络的误差函数  $err = \sum \|d - Y\|_2$  为二次函数, 如图 3 所示, 是一个一般二次函数及其导数的示意图, 通过二次函数示意图解释本文提出的自适应学习

率的特性. 本文提出基于权值变化的自适应学习率为  $\beta(t) = \frac{2}{1+e^{-|t| \times 10^n}} - 1$ , 其中  $t$  为权值的变化, 即误差曲面函数的梯度, 也就是如图 3 中的二次函数切线方程的斜率  $k$ , 于是  $t = k$ ,  $k = y/x$ ; 图中三点切线斜率分别为  $k_1, k_2, k_0$ , 从图中可以看到在点  $(x_0, y_0)$  处的导数的绝对值  $|k_0| < |k_1|$ ,  $|k_0| < |k_2|$ , 而学习率函数  $\beta(t) = \frac{2}{1+e^{-|t| \times 10^n}} - 1$ , 与  $t$  呈现正比例, 即与二次函数的斜率  $k$  呈现成比例, 则  $k$  对应的学习率分别为  $\beta_0, \beta_1, \beta_2$ , 于是有  $\beta_0 < \beta_1$ ,  $\beta_0 < \beta_2$ ; 在误差函数曲线没有收敛到极小值期间, 其权值变化为  $k_1$  (或者  $k_2$ ), 在此期间, 随着训练进行,  $k_1$  朝  $k_0$  逼近, 于是  $\beta_1$  向  $\beta_0$  逼近, 即学习率  $\beta$  渐变小, 在收敛的极小值点  $(X_0, Y_0)$  时, 学习率减小到  $\beta_0$ , 学习率达到最小值. 综上所述: 随着训练的进行, 权值的变化率逐渐变小, 学习率也逐渐变小, 达到自适应的目的.

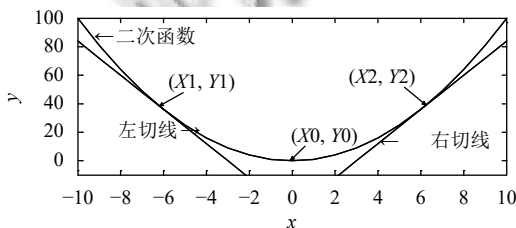


图 3 二次函数及其切线

对于现有固定学习率的神经网络, 学习率偏大, 容易产生震荡; 学习率偏小, 收敛速度慢, 网络拟合效果差, 不利于收敛; 本文提出的自适应学习率  $\beta$ , 根据权值变化自适应调整大小, 当权值变化大时, 此时学习率大; 当网络权值变化小, 学习率小 (如图 2); 在即将达到目标输出时, 误差接近极小值点, 误差曲面的梯度变化小, 即此时权值变化  $t$  较小, 从而学习率较小, 更有利于得到网络收敛, 对提高误差的精度, 具有显著的作用.

此外, 在每两个神经元之间, 其连接权值都有对应的学习率; 训练过程中, 每两个神经元的连接权值时刻在变化, 其对应的学习率也变化, 所以训练过程中, 产生数以万计的学习率, 以匹配权值的变化, 适应网络更新<sup>[10]</sup>.

针对现有的几种典型自适应学习率, 与本文的自适应学习率作对比:

1) 自适应全参数学习率 Adagrad<sup>[11-13]</sup> 是使学习率参数自适应变化, 把梯度的平方根作为学习率的分子, 训练前期梯度小, 则学习率大, 训练后期, 梯度叠加增

大, 学习率小; 由于累加的梯度平方根和越来越大, 学习率会逐渐变小, 最终趋于无限小, 严重影响网络收敛速度.

2) 牛顿法, 用 Hessian 矩阵替代学习率, 并结合梯度下降法, 虽然可得最优解, 但要存储和计算 Hessian 矩阵, 增大计算复杂度<sup>[14,15]</sup>.

3) 本文提出基于权值变化的自适应学习率, 利用参数  $10^n$  调整学习率对权值变化的敏感度, 以至于不存在如 Adagrad 算法的学习率趋于无限小的弊端; 本文的自适应学习率, 只需把权值更新过程中权值的导数用于学习率中, 计算的复杂度远低于牛顿法<sup>[16,17]</sup>.

### 2.3 权值更新

采用误差反向传播方式更新权值, 使误差  $e$  更快的取得极小值.

由式 (3)~(式 6)、式 (8) 得误差偏导为:

$$\frac{\partial \text{err}}{\partial W_j^L} = f(S_j^L) \times (1 - f(S_j^L)) \times (d - Y) \times X_n \quad (9)$$

对于基于自适应学习率的网络权值更新, 依梯度下降法得权值更新为:

$$W_n^L = W_j^L - \beta \times t \quad (10)$$

正态分布模型的权值偏导为:

$$\frac{\partial g}{\partial W_j^L} = \frac{\partial g}{\partial \text{err}} \times t \quad (11)$$

由式 (7)、式 (9)、式 (11) 可得:

$$\frac{\partial g}{\partial W_j^L} = \frac{-1}{\delta \times \sqrt{2\pi}} \times e^{-\left(\frac{\text{err}}{\delta}\right)^2} \times \frac{\text{err}}{\delta^2} \times (d - y) \times f(S_j^L) \times (1 - f(S_j^L)) \times X_n \quad (12)$$

式 (12) 可以看出, 网络训练后期, 误差  $\text{err}$  趋于极小的值, 此时权值变化不明显, 收敛速度慢; 为提高收敛速度, 提出解决方法为:

$$\text{err} = \text{sgn}(\text{err}) \times e^{|\text{err}|} \quad (13)$$

利用式 (13) 左边的  $\text{err}$  代替原来的误差  $\text{err}$ , 其中  $\text{sgn}(\text{err})$  为符号函数, 定义为:

$$\text{sgn}(\text{err}) = \begin{cases} 1, & \text{err} \geq 0 \\ -1, & \text{err} < 0 \end{cases} \quad (14)$$

对于正态分布模型, 依梯度上升法得权值更新:

$$W_n^L = W_j^L + \alpha \times t \quad (15)$$

其中  $W_n$  为  $W$  更新后的权值.

对于传统模型,依梯度下降法得权值更新为:

$$\begin{aligned} Wn_j^L &= W_j^L - \alpha \times t = W_j^L - \alpha \times \frac{\partial err}{\partial W_j^L} \\ &= W_j^L - \alpha \times \frac{\partial err}{\partial f(S_j^L)} \times \frac{\partial f(S_j^L)}{\partial S_j^L} \times \frac{\partial S_j^L}{\partial W_j^L} \\ &= W_j^L - \alpha \times (d-y) \times f(S_j^L) \times (1-f(S_j^L)) \times x_n \end{aligned} \quad (16)$$

正态分布模型与自适应学习率结合,依梯度上升法得权值更新为:

$$Wn_j^L = W_j^L + \beta \times \frac{\partial g}{\partial W_j^L} \quad (17)$$

结合式(7)、式(13)、式(14)、式(17)得权值更新为:

$$\begin{aligned} Wn_j^L &= W_j^L + \left( \frac{2}{1+e^{-|t| \times 10^m}} - 1 \right) \times \frac{1}{\delta^2} \times \frac{-1}{\delta \times \sqrt{2\pi}} \\ &\quad \times e^{-\left(\frac{err}{\delta}\right)^2} \times \text{sgn}(err) \times e^{|err|} \\ &\quad \times f(S_j^L) \times (d-y) \times (1-f(S_j^L)) \times X_n \end{aligned} \quad (18)$$

式(18)为结合梯度上升、正态分布模型、自适应学习率的权值更新方式,与传统权值更新方式(式(16))相比,改进后权值变化系数为:

$$\left( \frac{2}{1+e^{-|t| \times 10^m}} - 1 \right) \times \text{sgn}(err) \times \frac{1}{\delta^3 \times \sqrt{2\pi}} \times e^{|err| - \left(\frac{err}{\delta}\right)^2} \quad (19)$$

令  $\delta^2 = 2$ , 则  $e^{|err| - \left(\frac{err}{\delta}\right)^2} = e^{\frac{1-(|err|-1)^2}{2}}$ .

网络的参数经过归一化后,训练过程满足:  $0 \leq |err| \leq 1$ , 可得在网络训练初期,误差|err|较大,则  $e^{\frac{1-(|err|-1)^2}{2}} \in [1, e]$ 较大,权值更新  $Wn_j^L - W_j^L$ 较大,从而加快网络收敛速度,提高网络训练效率;训练后期,误差变小,该系数趋于1,对训练无明显影响。

### 3 实验验证

采用经典 XOR 问题,验证改进 BP 网络;标准 XOR 问题与验证 XOR 问题如表 1。

表 1 XOR 问题

标准输入 X	标准输出 Y	验证输入 X	验证输出 Y
0	0	0.1	0.2
0	1	0.15	0.9
1	0	1.1	0.01
1	1	0.88	1.03

依据实验得出结果

先用标准 XOR 问题对神经网络训练,再用接近标准 XOR 输入对神经网络验证,比较验证输出与标准输出,判断优劣。

#### 3.1 带正态分布的固定学习率模型与传统模型对比

设定学习率: 0.5, 误差限默认: 0.000 001, 迭代次数默认 10 000 次,得基于正态分布模型的 BP 网络和传统 BP 网络误差曲线,如图 4。

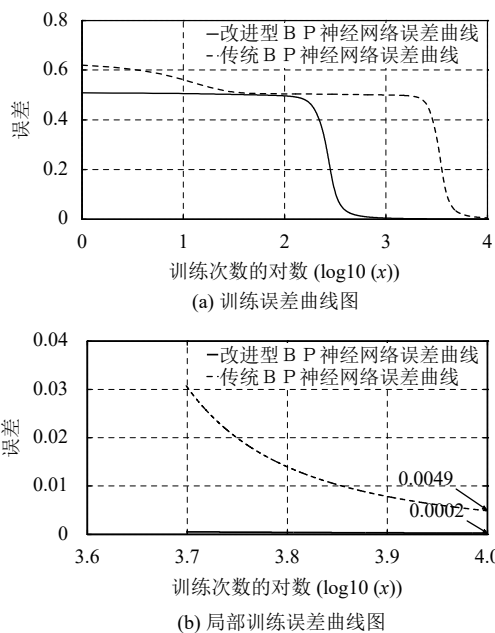


图 4 正态分布模型与传统模型的误差对比

XOR 异或问题的验证输出为表 2。

表 2 验证输出 Y

验证输入 X	传统模型输出 Y	改进模型输出 Y	标准输出 Y
0.1 0.2	0.1496	0.0122	0
0.15 0.9	0.8988	0.9499	1
1.1 0.01	0.9748	0.9934	1
0.88 1.03	0.0646	0.0161	0

从图 4、表 2 可以看出,基于正态分布模型改进后的网络,其误差是传统模型的 1/25,误差明显降低,且验证结果更接近于标准输出。

分别比较不同学习率和不同训练次数之间的误差,如表 3。

从大量实验可以看出,基于正态分布模型的 BP 网络与传统 BP 网络模型相比,具有更小的误差或更快的迭代速度。(带\*为改进模型实验误差,带\*\*为传统模型实验误差)。



表3 不同学习率与训练次数的误差对比

训练次数 (次)	学习率			
	0.05	0.1	0.3	0.5
1000	0.4960*	0.4342*	0.0107*	0.0041*
	0.5045**	0.5033**	0.5010**	0.4998**
5000	0.0167*	0.0041*	0.0009*	0.0005*
	0.5014**	0.4998**	0.3725**	0.0308**
10 000	0.0049*	0.0015*	0.0004*	0.0002*
	0.4998**	0.4919**	0.0162**	0.0049**
20 000	0.0015*	0.0007*	0.0002*	0.0001**
	0.4919**	0.0958**	0.0035**	0.0015**

3.2 自适应学习率模型与固定学习率模型对比

自适应学习率的倾斜参数  $n$  设为 3. 以权值变化作为自适应学习率变化依据, 采用梯度下降法更新权重. 得隐含层自适应学习率变化曲线, 图 5 所示.

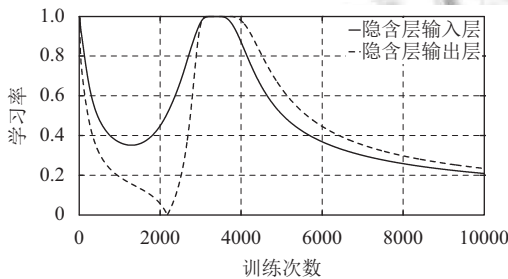


图5 适应性学习率的变化

可以看出, 学习率是随权值的变化而自适应变化, 每一轮迭代后, 权值变化不同, 导致学习率不同; 训练后期, 权值变化减小, 学习率减小, 自适应学习率相应减小.

对于固定学习率, 采用自适应学习率的算术平均值: 0.4567, 将改进型自适应学习率与固定学习率训练结果作对比, 如图 6.

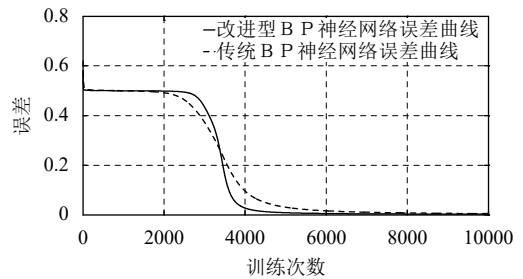
XOR 异或问题的验证输出为表 4.

可以看出, 自适应学习率模型的误差为固定学习率模型的 1/2.2, 并且验证结果更接近标准 XOR 异或问题.

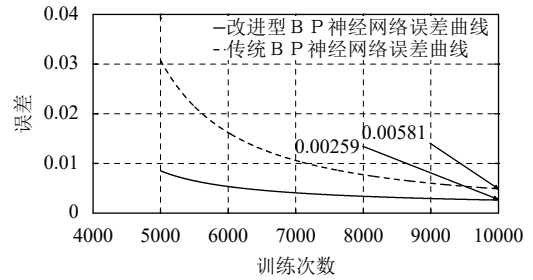
3.3 自适应学习率的正态分布模型与固定学习率模型对比

倾斜系数为 3, 可得适应性学习率的变化与训练次数之间的关系为图 7.

固定学习率采用自适应学习率的算术平均值: 0.1459, 与带自适应学习率和正态分布模型的 BP 神经网络对比, 如图 8.



(a) 训练误差曲线图



(b) 局部训练误差曲线图

图6 误差对比

表4 验证输出 Y

验证输入 X	传统模型输出 Y	改进模型输出 Y	标准输出 Y
0.1 0.2	0.1596	0.0480	0
0.15 0.9	0.8924	0.8317	1
1.1 0.01	0.9428	0.9729	1
0.88 1.03	0.0704	0.0498	0

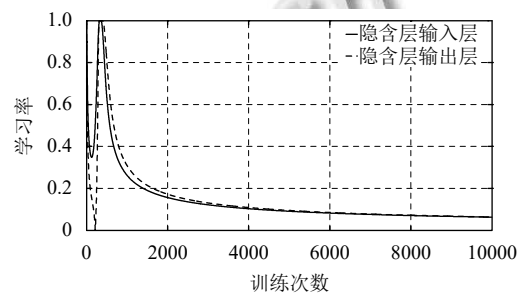


图7 适应性学习率的变化

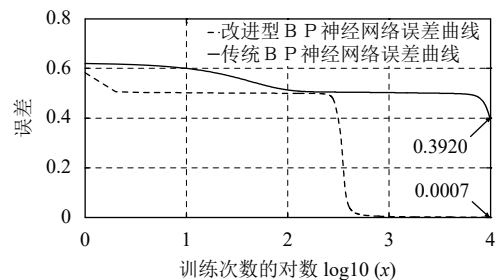


图8 误差对比

XOR 异或问题的验证输出为表 5。

表 5 验证输出  $Y$

验证输入 $X$	传统模型输出 $Y$	改进模型输出 $Y$	标准输出 $Y$
0.1 0.2	0.0992	0.0252	0
0.15 0.9	0.9286	0.8784	1
1.1 0.01	0.9681	0.9864	1
0.88 1.03	0.0401	0.0261	0

从图 8、表 5 可以看出,改进 BP 神经网络的误差是传统模型的 1/55,改进的 BP 神经网络性能明显优于传统模型。

#### 4 结束语

本文提出基于权值变化的自适应学习率、结合梯度上升法的正态分布模型,提升 BP 神经网络的运算效率;理论分析了提高收敛速度、降低误差的原理,通过仿真结果表明,改进后的 BP 神经网络在提高收敛速度、降低误差方面具有更好的成效。

#### 参考文献

- Li LS, Gan SJ, Yin XD. Feedback recurrent neural network-based embedded vector and its application in topic model. *EURASIP Journal on Embedded Systems*, 2017, (2017): 5. [doi: 10.1186/s13639-016-0038-6]
- 李新叶, 黄腾. 基于多尺度跃层卷积神经网络的精细车型识别. *科学技术与工程*, 2017, 17(11): 246–249. [doi: 10.3969/j.issn.1671-1815.2017.11.041]
- He W, Chen YH, Yin Z. Adaptive neural network control of an uncertain robot with full-state constraints. *IEEE Transactions on Cybernetics*, 2016, 46(3): 620–629. [doi: 10.1109/TCYB.2015.2411285]
- Li XW, Cho SJ, Kim ST. Combined use of BP neural network and computational integral imaging reconstruction for optical multiple-image security. *Optics Communications*, 2014, (315): 147–158. [doi: 10.1016/j.optcom.2013.11.003]
- 杨志浩, 李治平. 基于 BP 神经网络的底水油藏控水压裂选段新方法. *地质与勘探*, 2017, 53(4): 818–824.
- 李奎, 李晓倍, 郑淑梅, 等. 基于 BP 神经网络的交流接触器剩余电寿命预测. *电工技术学报*, 2017, 32(15): 120–127.
- Dong RL, Zhao GY. The use of artificial neural network for modeling in vitro rumen methane production using the CNCPS carbohydrate fractions as dietary variables. *Livestock Science*, 2014, (162): 159–167. [doi: 10.1016/j.livsci.2013.12.033]
- Yu F, Xu XZ. A short-term load forecasting model of natural gas based on optimized genetic algorithm and improved BP neural network. *Applied Energy*, 2014, (134): 102–113. [doi: 10.1016/j.apenergy.2014.07.104]
- Ahmed R, El Sayed M, Gadsden SA, et al. Automotive internal-combustion-engine fault detection and classification using artificial neural network techniques. *IEEE Transactions on Vehicular Technology*, 2015, 64(1): 21–33. [doi: 10.1109/TVT.2014.2317736]
- Jia WK, Zhao DA, Shen T, et al. An optimized classification algorithm by BP neural network based on PLS and HCA. *Applied Intelligence*, 2015, 43(1): 176–191. [doi: 10.1007/s10489-014-0618-x]
- 李雪芝, 周建平, 许燕, 等. 基于 L-M 算法的 BP 神经网络预测短电弧加工表面质量模型. *燕山大学学报*, 2016, 40(4): 296–300, 318.
- 赵一鹏, 丁云峰, 姚恺丰. BP 神经网络误差修正的电力物资时间序列预测. *计算机系统应用*, 2017, 26(10): 196–200. [doi: 10.15888/j.cnki.csa.006011]
- 贾楠, 胡红萍, 白艳萍. 基于 BP 神经网络的人口预测. *山东理工大学学报(自然科学版)*, 2011, 25(3): 22–24.
- Hernández-Pajares M, Juan JM, Sanz J. Neural network modeling of the ionospheric electron content at global scale using GPS data. *Radio Science*, 1997, 32(3): 1081–1089. [doi: 10.1029/97RS00431]
- 潘庆先, 董红斌, 韩启龙, 等. 一种基于 BP 神经网络的属性重要性计算方法. *中国科学技术大学学报*, 2017, 47(1): 18–25.
- 胡燕祝, 李雷远. Kalman 滤波-BP 神经网络在执行机构自主定位中的应用. *北京邮电大学学报*, 2016, 39(6): 110–115.
- 崔丽辉, 赵安兴, 宁方正. 基于 EMD 和 BP 神经网络的雷达体特征信号检测算法. *计算机系统应用*, 2017, 26(8): 217–222. [doi: 10.15888/j.cnki.csa.005920]