

# 一种学习速率自适应的可编程片上学习 BP 神经网络电路系统的设计

卢 纯, 石秉学, 陈 卢

(清华大学微电子学研究所, 北京 100084)

**摘 要:** 设计了一种学习速率自适应的可编程片上学习 BP 神经网络电路系统. 整个系统由前向网络、误差反传网络两部分组成. 提出了一种新型的可编程 S 型函数及其导数发生器电路. 它不仅产生 S 型函数, 完成非线性 I-V 转换; 还利用前向差分法, 产生 S 型函数的导数. 这两种函数不仅与理想函数的拟合程度很好, 而且易实现对阈值和增益因子的编程. 为提高 BP 神经网络片上学习的收敛速度, 还提出了学习速率自适应电路. 本文采用标准 112Lm CMOS 工艺的模型参数, 对整个系统进行了  $\sin(x)$  函数拟合等模拟实验, 验证了该片上学习 BP 神经网络的优越性能.

**关键词:** 神经网络; CMOS 模拟集成电路; 可编程; 自适应

**中图分类号:** TP18      **文献标识码:** A      **文章编号:** 0372-2112 (2001) 05-0702-03

## A Circuit System Design of Programmable BP On-Chip Learning Neural Network with Learning Rate Adaptation

LU Chun, SHI Bingxue, CHEN Lu

(Institute of Microelectronics, Tsinghua University, Beijing 100084, China)

**Abstract:** A circuit system of programmable BP on-chip learning neural network with learning rate adaptation is designed. The whole system comprises feedforward network and error back-propagation network. A novel programmable generator of sigmoidal function and its derivation is proposed. Its outputs include the sigmoidal function to realize I-V nonlinear transfer and its derivative using the forward differential method. Both functions fit well with the ideal functions. Moreover, the threshold and the gain factor can be easily programmable. Learning rate adaptation circuit is also presented to accelerate the convergent speed. Using a standard 112Lm CMOS process, experiments such as  $\sin(x)$  function fitness are done to the whole system. These experiments verify the superior performance of this on-chip learning BP neural network.

**Key words:** neural networks; CMOS analogue integrated circuits; programmable; adaptation

### 1 引言

神经网络的硬件实现是当今神经网络领域重要的、不可或缺的重要组成部分, 特别适用于以下的应用场合: (1) 体积小; (2) 重量轻; (3) 适应环境变化; (4) 速度快; (5) 规模大. 神经网络的硬件实现中, 有监督学习类型分为三种: 片下学习 (off-chip learning), 芯片循环学习 (chip-in-the-loop learning) 和片上学习 (on-chip learning). 片下学习在芯片外完成所有的学习过程. 学习完毕后, 再将权重加载到芯片上. 因为用于训练的主机和神经网络芯片的精度不同, 所以加载时, 要对权重进行四舍五入. 权重的四舍五入和芯片的非理想特性使权重加载后的系统实际性能可能会很差. 芯片循环学习中, 芯片参与前向运算, 而权重的迭代在片外完成, 这样就补偿了特定芯片的非理想特性. 但片内与片外数据间的传输成为影响学习速度的瓶颈. 而且由于不能完全脱离主机, 其应用范围也受到一定限制. 由于片下学习和芯片循环学习存在以上的缺点, 现在很多学者已开始对片上学习<sup>[1]</sup>进行深入研究. 片上学习中, 前向运算和权重的更新都在芯片上完成. 它的优势在于以下 3 点: (1) 快速: 片上学习不仅将神经网络本身, 而且将学习算法映

射成硬件, 全方位地利用了神经网络及其算法的并行特性, 从而大大缩短了权重迭代的时间. (2) 完全脱机: 整个训练过程无需主机, 使神经网络更具独立性, 更能适应环境变化和实时处理的要求, 也更利于充分发挥神经网络硬件实现体积小、重量轻的优点. (3) 对电路非理想特性的补偿作用: 将神经网络的前向运算和权重迭代统一用硬件实现, 几乎可以补偿所有由电路非理想特性引起的偏差.

片上学习有两种最常用的算法: 误差反传算法 (Back-Propagation algorithm, 简称 BP 算法) 和权重扰动算法<sup>[2]</sup> (Weight Perturbation algorithm, 简称 WP 算法). BP 算法在神经网络的研究中有较长的历史, 它最早是由 R. J. Werbos 在 1974 年提出的, Rumelhart 等于 1985 年发展了 BP 算法, 使其逐渐受到人们的重视. BP 的优点在于所有的权重都是并行调整的, 而且权重的调整只需一些局部的信息. WP 算法由 M. Jabri 等人于 1992 年提出, 它的特点在于只需前向运算过程, 而无需反向运算; 缺点在于权重的调整是串行的. 虽然后来又有人提出了半并行 WP 算法, 但为了充分发挥神经网络硬件实现并行度高, 速度快的优点, 本文采用 BP 算法作为片上学习的算法. 通

过对神经元阈值和增益因子的可编程设计,该神经网络将可适用到广泛的场合;而学习速率的自适应使该网络在保证收敛几率的基础上使收敛速度大大提高.

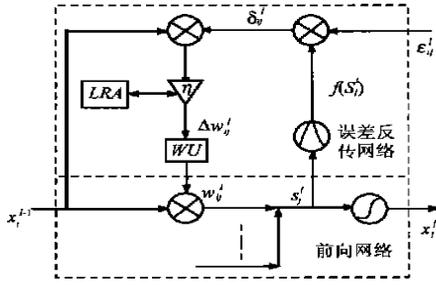


图 1 片上学习 BP 神经网络整体电路结构框图

## 2 电路系统结构

BP 网络包括输入层、隐含层和输出层. 各个节点的转移函数一般都采用 S 型函数, 它的表达式为

$$f(s) = 1 / (1 + e^{-As}) \quad (1)$$

其中 A 为增益因子, s 为某结点所有输入的加权和. 设 L (LE 1) 层 BP 网络的训练样本数为 R, 第 1 (l = 1, 2, ..., L) 层第 j 个神经元的阈值为  $H_j^l$ , 第 1 (l = 1, 2, ..., L) 层的第 i (0F i < n) 个神经元与它之间的权重为  $w_{ij}^l$ . 为方便起见, 通过在第 l-1 层中加入输入恒为 1 的第 n 个神经元, 将阈值写入权重中, 即  $x_n^{l-1} = 1, H_j^l = w_{nj}^l$ . 对于训练样本 r (r = 1, 2, ..., R) 来说, 设第 l-1 层的第 i 个神经元的输出为  $x_{i,r}^{l-1}$ ; 第 l 层第 j 个神经元的输出为  $x_{j,r}^l, l = L$  时, 训练信号为  $t_{j,r}$ ; 第 l-1 层的所有神经元的输出到第 l 层第 j 个神经元的加权和为  $s_{j,r}^l$ , 则网络前向运算由式(2)表示:

$$x_{j,r}^l(k) = f(s_{j,r}^l(k)) = f\left(\sum_{i=0}^n w_{ij}^l(k) x_{i,r}^{l-1}(k)\right) \quad (2)$$

在误差反传及权重迭代过程中, 定义权重误差  $D_{ij,r}^l(k) = f'(s_{j,r}^l(k)) E_{j,r}^l(k)$ , 其中神经元误差  $E_{j,r}^l(k) = \begin{cases} t_{j,r} - x_{j,r}^l(k), & l = L \\ E_j w_{ij}^{l+1}(k) D_{ij,r}^{l+1}(k), & 1 \leq l < L \end{cases}$ , 则权重变化量  $\Delta w_{ij}^l(k+1) = \sum_{r=1}^R D_{ij,r}^l(k) x_{i,r}^{l-1}(k)$ , 误差反传运算由式(3)表示:

$$w_{ij}^l(k+1) = w_{ij}^l(k) + G_j^l(k+1) \sum_{r=1}^R D_{ij,r}^l(k) x_{i,r}^{l-1}(k) \quad (3)$$

式(3)中的  $G_j^l$  为学习速率. 它对收敛速度的影响很大. 本文采用局部的、自适应的学习速率  $G_j^l$ .  $G_j^l$  自适应的过程中, 遵循以下原则: 若相邻两次迭代中, 权值的变化  $\Delta w_{ij}^l$  的符号相同, 则增大  $G_j^l$ ; 若相邻两次迭代中, 权值的变化  $\Delta w_{ij}^l$  的符号相反, 则减小  $G_j^l$ . 也就是说, 设  $d_{ij}^l(k+1)$  为  $\Delta w_{ij}^l(k+1)$  的符号, 若  $d_{ij}^l(k+1) = d_{ij}^l(k)$ , 则,  $G_j^l(k+1) = G_j^l(k) \left( G_j^{\max} / G_j^l(k) \right)^c$  (4) 否则,  $G_j^l(k+1) = G_j^l(k) \left( G_j^{\min} / G_j^l(k) \right)^c$  (5)

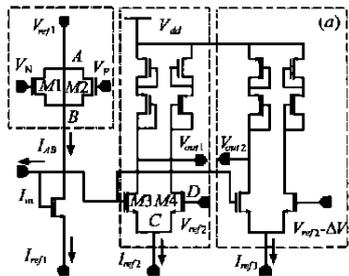
根据以上算法, 可将电路系统粗略地分为前向网络、误差反传网络两部分. 训练样本数 R 为 1 时的电路结构如图 1 所示, 图中标注的各变量为前面算法中相应变量略去脚标 r 所得. 前向网络主要包括突触和神经元两种基本处理单元. 突触实现加权求和功能, 采用了简单、占用面积小的传统的四象限 Gilbert 乘法器. 神经元主要完成非线性 I-V 转换. 误差反传网络包括 S 型函数的导数发生器, 权重迭代单元 WU (Weight Updating unit) 和学习速率自适应电路 LRA (Learning Rate Adaptation circuit) 等. 片上学习中, 误差反传网络部分的规模远远大于前向神经网络本身. 对误差反传网络部分电路设计的研究也远没有前向网络成熟. 因此, 误差反传网络部分成为本文讨论的重点, 提出了一种新型的可编程 S 型函数及其导数的发生器和 LRA 电路.

## 3 单元电路设计

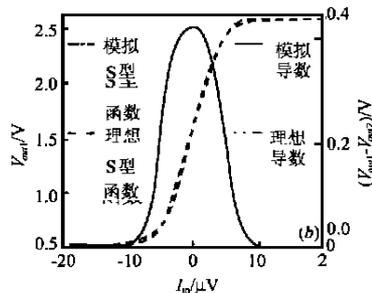
### 3.1 可编程 S 型激活函数及其导数的发生器

可编程 S 型激活函数及其导数的发生器电路<sup>[4]</sup>如图 2 (a) 所示.  $V_{dd}$  为 3.13V 电压源. 虚线框中 N 管 M1 和 P 管 M2 组成线性可调电阻,  $V_N, V_P$  为其控制信号.  $V_{ref1}$  为固定电平, 它的选择保证 M1 和 M2 处于线性状态. 每一个点划框中, 两个 N 管组成的差分输入放大器和 4 个 P 管组成的负载一起实现非线性 I2V 变换. 差分输入的一端与 B 点相接, 另一端为固定电平  $V_{ref2}$  或  $V_{ref2} - \Delta V$ , 这里  $\Delta V$  为固定小电压.  $I_{ref1}, I_{ref2}$  为恒流源电流, 它们的参考方向如图 2 (a) 箭头中所示. 由文献[4]可知,  $V_{out}$  输出 S 型激活函数, 且与理想 S 型激活函数的误差不超过 3%, 如图 2 (b) 中点划线和虚线所示.

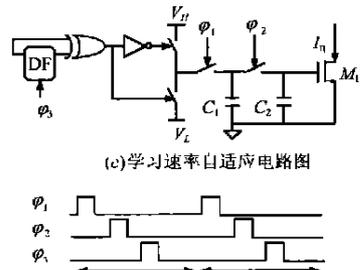
设电路产生的激活函数表达为  $V_{out} = V_{out}(I_{in})$ . 运用前向差分法, 转移函数的导数  $V_{der iv}$  可由  $V_{out1}$  减去  $V_{out2}$  得到:



(a) 神经元电路图



(b) 模拟 S 型函数及其导数与理想函数间的比较



(c) 学习速率自适应电路

第 k 次迭代 第 k+1 次迭代

图 2 S 型函数及其导数发生器与学习速率自适应电路

$$V_{out}(I_{in}) = V_{out}(V_d) \# V_d(I_{in})$$

$$\mu = \frac{V_{out}(V_B - V_{ref2} + S V) - V_{out}(V_B - V_{ref2})}{S V} \# R_{AB} \quad (6)$$

$$V_{d \text{ deriv}}(I_{in}) S \frac{S V}{R_{AB}} \# V_{out}(I_{in}) \mu - (V_{out}(V_B - (V_{ref2} - S V)))$$

$$- V_{out}(V_B - V_{ref2}) = V_{out1} - V_{out2} \quad (7)$$

采用标准 112Lm CMOS 工艺的模型参数, 对  $(V_{out1} - V_{out2})$  进行的 HSPICE 模拟结果如图 2(b) 中实线所示。图 2(b) 中点线为理想导数。它们间的相对误差不超过 5%。

片上学习神经网络的优势之一就是它能适应未知的和变化的环境。因此, 电路良好的编程性能是很重要的。图 2(a) 所示的新型神经元可实现阈值和增益因子的编程。改变  $I_{ref1}$  可以调节阈值。若  $I_{ref1}$  增大, 转移曲线及其导数曲线向左平移; 反之, 则向右平移。另外, 改变  $V_N$ 、 $V_P$  的值起到调节增益因子大小的作用。图 2(a) 中 M1、M2 工作在线性区, 取  $B_1 = B_2$ , 则

$$R_{AB} = \frac{V_A - V_B}{I_1 + I_2} = \frac{1}{B[(V_N - V_P) - (V_{T1} + |V_{T2}|)]} \quad (8)$$

由式(8)可知,  $(V_N - V_P)$  越大,  $R_{AB}$  越小,  $V_B$  随输入电流变化的斜率越小,  $V_{out1}$  上升越缓慢, 式(1)中增益因子 A 越小。

### 3.1.2 学习速率自适应电路

如图 2(c) 所示的学习速率自适应电路是式(4)和式(5)的硬件映射。  $V_s$  为权重变化的符号位。  $I_G$  为输出学习速率。 NMOS 管 M1 工作于弱反型区。参考电压  $V_H$ 、 $V_L$  ( $V_H > V_L$ ) 用来设置  $G^{max}$  和  $G^{min}$ 。  $U_1$ 、 $U_2$  和  $U_3$  是三相不重叠时钟, 如图 2(d) 所示。在第 k 次迭代时,  $U_3$  的下降沿使  $V_s(k)$  传输到 D 触发器的输出端。在第 k+1 次迭代时,  $V_s(k)$  与  $V_s(k+1)$  一起送到异或门的输入端。这样, 在  $U_1$  为高电平时, 若  $V_s(k) = V_s(k+1)$ , 则  $V_{C1}(k+1) = V_H$ , 否则  $V_{C1}(k+1) = V_L$ 。  $U_2$  为高电平时, 通过电容  $C_1$ 、 $C_2$  间的电荷共享, 实现学习速率的自适应。基本原理如下:

$$V_{C2}(k+1) = \frac{C_2}{C_1 + C_2} V_{C2}(k) + \frac{C_1}{C_1 + C_2} V_{C1}(k+1)$$

$$= V_{C2}(k) + \frac{C_1}{C_1 + C_2} (V_{C1}(k+1) - V_{C2}(k)) \quad (9)$$

设  $C = \frac{C_1}{C_1 + C_2}$ , M1 工作与弱反型区, 则

$$I_G(k+1) = I_S e^{V_{C2}(k+1)/nU_T} = I_G(k) \left( \frac{I_{C1}(k+1)}{I_G(k)} \right)^C \quad (10)$$

其中,  $I_{C1}(k+1) = I_S e^{V_{C1}(k+1)/nU_T}$ 。设  $I_G^{max} = I_S e^{V_H/nU_T}$ ,  $I_G^{min} = I_S e^{V_L/nU_T}$  那么, 若  $V_s(k) = V_s(k+1)$ ,  $I_G(k+1) = I_G(k) \left( \frac{I_G^{max}}{I_G(k)} \right)^C$ ; 否则,  $I_G(k+1) = I_G(k) \left( \frac{I_G^{min}}{I_G(k)} \right)^C$ 。对照式(4)和式(5)可知, 图 2(c) 所示电路完成了学习速率的自适应。

## 4 实验结果

函数拟合是检验通用神经网络功能的基准之一。本文对  $\sin(x)$  函数的完整周期进行了拟合。实验所采用的网络结构为输入层 1 个神经元、隐含层 5 个神经元、输出层 1 个神经元。实验结果如图 3 所示。实心圆点代表训练集中的 64 个训练元素。空心圆点为网络的实际输出。实验中发现隐含层神经元的个数太少或太多, 都易使学习陷入极小点。对于  $\sin(x)$  函

数完整周期拟合问题来说, 隐含层神经元的个数取 4-6 个比较合适。从图 3 可以看出, 采用 1-5-1 的拓扑结构, 本片上学习系统能很好地拟合  $\sin(x)$  函数。

## 5 结论

近年来, 片上学习以其快速、完全脱机及对芯片非理想特性的补偿等优良特性受到了神经网络研究领域的重视。本文设计了一种学习速率自适应的可编程片上学习 BP 模

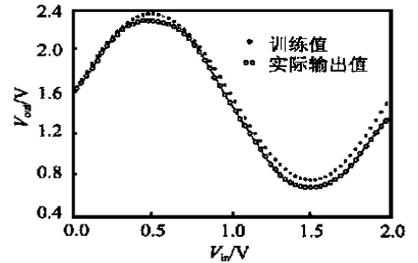


图 3  $\sin(x)$  函数拟合实验结果

拟神经网络电路系统; 提出了一种新型的可编程 S 型函数及其导数发生器; 设计了学习速率自适应电路。S 型函数及其导数发生器利用前向差分原理, 同时产生 S 型函数及其导数。模拟结果表明, 该发生器产生的激活函数与理想 S 型函数的相对误差不超过 3%; 产生的导数与理想导数的相对误差不超过 5%。通过对其阈值和增益因子的编程, 该发生器还可根据不同的要求, 设置不同的特性, 使神经网络的片上学习更具适应环境的能力。在学习速率自适应电路的设计中, 兼顾了三方面的因素: (1) 每个权重可根据自身不同的局部特性进行学习速率的调节; (2) 电路易于 VLSI 实现; (3) 可被推广应用到其它算法的硬件实现中。采用标准 112Lm CMOS 工艺的模型参数, 对整个系统进行的  $\sin(x)$  函数拟合实验有效地验证了该片上学习 BP 神经网络的优越性能。

### 参考文献:

[ 1 ] M. Valle, D. D. Caviglia and G. M. Bisio. An analog VLSI neural net2 work with on2chip back propagation learning [ J ]. Analog Integrated Circuits and Signal Processing 1996, 9: 231- 245.

[ 2 ] M. Jabri and B. Flower. Weight perturbation: an optimal architecture and learning technique for analog VLSI feedforward and recurrent mul2 tiplayer networks [ J ]. IEEE Trans. On Neural Network, 1992, 3 ( 1 ): 154- 157.

[ 3 ] B. K. Dolenko and H. C. Card. Tolerance to analog hardware of on2chip learning in backpropagation networks [ J ]. IEEE Trans. On Neural Net2 works, 1995, 6 ( 5 ): 1045- 1052.

[ 4 ] Lu, C. and Shi, B. X., Circuit design of an adjustable neuron activation function and its derivative [ J ]. Electronics Letters, 2000, 36 ( 6 ): 553 - 555.

### 作者简介:



卢 纯 1975 年出生, 1997 年取得清华大学电子工程系微电子专业工学学士学位。同年获得直接攻读清华大学博士研究生资格。现为清华大学博士研究生, 从事人工神经网络的 VLSI 实现和模拟集成电路研究。