

多层感知器模型互反奇异性区域学习动态的理论分析

郭伟立¹, 魏海坤^{1†}, 赵军圣^{1,2}, 张侃健¹

(1. 东南大学 自动化学院 复杂工程系统测量与控制教育部重点实验室, 江苏 南京 210096;

2. 聊城大学 数学科学学院, 山东 聊城 252059)

摘要: 多层感知器神经网络(MLPs)的学习过程经常发生一些奇异性行为, 容易陷入平坦区, 这都和MLPs的参数空间中存在的奇异性区域有直接关系. 当MLPs的两个隐节点的权值接近互反时, 置换对称性会导致学习困难. 对MLPs的互反奇异性区域附近的学习动态进行分析. 本文首先得到了平均学习方程的解析表达式, 然后给出了互反奇异性区域附近的理论学习轨迹, 并通过数值方法得到了其附近的实际学习轨迹. 通过仿真实验, 分别观察了MLPs的平均学习动态, 批处理学习动态和在线学习动态, 并进行了比较分析.

关键词: 多层感知器; 神经网络; 学习动态; 奇异性; 互反

中图分类号: TP273 文献标识码: A

Theoretical analysis of learning dynamics near the opposite singularities in multilayer perceptrons

GUO Wei-li¹, WEI Hai-kun^{1†}, ZHAO Jun-sheng^{1,2}, ZHANG Kan-jian¹

(1. Key Laboratory of Measurement and Control of CSE, Ministry of Education, School of Automation, Southeast University, Nanjing Jiangsu 210096, China;

2. School of Mathematics Science, Liaocheng University, Liaocheng Shandong 252059, China)

Abstract: Owing to the existence of singularities in the parameter space, multilayer perceptrons (MLPs) may behave extremely slowly in learning or even be trapped in plateaus. When weights of two hidden units are nearly mutually opposite, the learning process will encounter difficulties because of the permutation symmetry. We investigate the learning dynamics of MLPs near opposite singularities, and derive the analytical expressions for averaged learning equations. Then, we obtain the theoretical learning trajectories near the opposite singularities. Furthermore, real learning trajectories near the opposite singularities are also calculated by using numerical methods. In simulations, we study the averaged learning dynamics, the batch mode learning dynamics and the online learning dynamics, respectively.

Key words: multilayer perceptrons; neural networks; learning dynamics; singularity; opposite

1 引言(Introduction)

笔者称学习机器是正则的: 如果其Fisher信息阵是可逆矩阵, 反之笔者称学习机器是奇异的. 研究表明包括前馈神经网络在内的几乎所有的学习机器都是奇异的^[1]. 在分层系统的奇异性区域中, 由于Fisher信息阵退化, 对梯度下降法而言负梯度方向不再是下降速度最快的方向, 学习过程中Cramér-Rao定理不再成立^[2-4]. 此时学习过程就会发生一些奇异性行为, 如容易陷入局部极小点, 有时学习速度变得很慢, 发生平坦区现象(如图1)^[5-10]. 受奇异性区域的影响, 赤池信息准则(AIC)、贝叶斯信息准则(BIC)、最小描述长度(MDL)等模型选择准则也不再有效. 文献[11-12]研究了奇异性区域在Bayes推断中的影响, 并提出了

一种应用更广泛BIC(WBIC)准则.

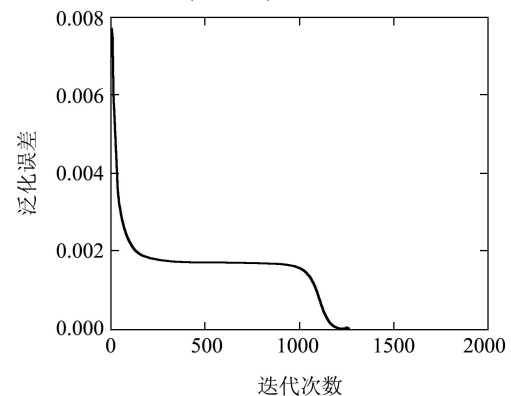


图 1 学习过程发生平坦区现象

Fig. 1 Plateau phenomenon occurred in the learning process

收稿日期: 2013-06-18; 录用日期: 2013-10-10.

[†]通信作者. E-mail: hkwei@seu.edu.cn.

基金项目: 国家自然科学基金重大项目资助项目(11190015); 高等学校博士学科点专项科研基金资助项目(20100092110020); 国家自然科学基金资助项目(61374006).

对典型的前馈神经网络

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{i=1}^k w_i \phi(\mathbf{x}, \mathbf{J}_i), \quad (1)$$

其中: \mathbf{J}_i 为输入层到隐层的权值, w_i 为隐层到输出层的权值. 当 w_i (或 $\|\mathbf{J}_i\|$) = 0时, $w_i \phi(\mathbf{x}, \mathbf{J}_i) = 0$, 此时无论 $\|\mathbf{J}_i\|$ (或 w_i)的取值为何, 对 $f(\mathbf{x}, \boldsymbol{\theta})$ 都不会造成影响, 则在神经网络的学习过程中就无法辨识 $\|\mathbf{J}_i\|$ (或 w_i). 另外, 当神经网络的参数 $\mathbf{J}_i = \pm \mathbf{J}_j$ 时, $w_i \phi(\mathbf{x}, \mathbf{J}_i) + w_j \phi(\mathbf{x}, \mathbf{J}_j) = (w_i \pm w_j) \phi(\mathbf{x}, \mathbf{J}_i)$, 这两个隐节点可以看作权值为 $\mathbf{J} = \mathbf{J}_i = \pm \mathbf{J}_j$ 和 $w = w_i \pm w_j$ 的一个隐节点, 此时系统只能辨识出两个输出权值的和(或差) $w_i \pm w_j$, 但不能分别分辨出 w_i 和 w_j . 所以前馈神经网络中存在两类奇异性区域^[13-14]:

$$1) \quad w_i \|\mathbf{J}_i\| = 0, \quad (2)$$

$$2) \quad \mathbf{J}_i = \pm \mathbf{J}_j. \quad (3)$$

多层感知器神经网络(multilayer perceptrons, MLPs)和径向基函数(radial basis function, RBF)网络都是典型的前馈神经网络, 在实际中都获得了广泛的应用. 研究表明RBF神经网络中只存在重合奇异性($\mathbf{J}_i = \mathbf{J}_j$)和消去奇异性($w_i = 0$)这两类奇异性区域^[2]. 文献[15]对分层网络的重合奇异性区域($\mathbf{J}_i = \mathbf{J}_j$)附近的学习动态做了通用的数学分析, 得到了重合奇异性区域附近的通用学习轨迹. 使用文献[15]中的方法, 文献[16]得到了RBF网络的平均学习方程, 分析了重合奇异性区域附近的学习动态, 并提出了奇异性区域附近平坦区现象发生的机制. 文献[17]讨论了toy模型中重合奇异性区域附近的学习动态, 该toy模型由两个隐节点的学生模型学习一个隐节点的教师模型, 并且所有参数的维数设定为1. 以上研究成果都是以重合奇异性区域为主要研究对象. 而当MLPs的激活函数是双极性函数, 即 $\phi(\mathbf{x}, -\mathbf{J}) = -\phi(\mathbf{x}, \mathbf{J})$ 时, 根据前面的分析, MLPs的参数空间还存在着互反奇异性区域 $\mathcal{R} = \{\boldsymbol{\theta} | \mathbf{J}_i = -\mathbf{J}_j\}$.

在本文中, 主要考虑MLPs中参数空间的互反奇异性区域:

$$\mathcal{R} = \{\boldsymbol{\theta} | \mathbf{J}_i = -\mathbf{J}_j\}. \quad (4)$$

首先求得MLPs的平均学习方程的解析表达式. 应用坐标变换, 得到互反奇异性区域附近的渐近学习轨迹, 并得到实际的学习轨迹. 在仿真实验中分别研究MLPs的平均学习动态、批处理学习动态和在线学习动态, 并对仿真结果进行分析.

2 问题描述(Problem statement)

文献[16]中对RBF网络的研究结果表明研究具有两个隐节点的模型已经足够, 所以本文主要研究具有两个隐节点的MLPs:

$$f(\mathbf{x}, \boldsymbol{\theta}) = w_1 \phi(\mathbf{x}, \mathbf{J}_1) + w_2 \phi(\mathbf{x}, \mathbf{J}_2), \quad (5)$$

其中: $\phi(\cdot)$ 是双极性的激活函数, $\phi(\mathbf{x}, \mathbf{J}_i) = \phi(\mathbf{J}_i^T \mathbf{x})$ ($i = 1, 2$)是第 i 个隐节点的输出, $\boldsymbol{\theta} = \{w_1, w_2, \mathbf{J}_1, \mathbf{J}_2\}$ 表示所有的系统参数.

假设教师模型也表示为具有两个隐节点的MLPs:

$$y = f_0(\mathbf{x}) = f(\mathbf{x}, \boldsymbol{\theta}_0) + \varepsilon = w_1 \phi(\mathbf{x}, \mathbf{t}_1) + w_2 \phi(\mathbf{x}, \mathbf{t}_2) + \varepsilon, \quad (6)$$

其中: ε 是服从0均值高斯分布的附加噪声, 与训练样本 \mathbf{x} 不相关, $\boldsymbol{\theta}_0 = \{w_1, w_2, \mathbf{t}_1, \mathbf{t}_2\}$ 表示教师的所有参数. 当教师函数不能由MLPs表示时, $f(\mathbf{x}, \boldsymbol{\theta}_0)$ 假定为MLPs表示的最佳逼近.

不失一般性, 假设输入 \mathbf{x} 服从均值为0和方差为单位阵 \mathbf{I}_n 的高斯分布:

$$q(\mathbf{x}) = (\sqrt{2\pi})^{-n} \exp\left(-\frac{\|\mathbf{x}\|^2}{2}\right). \quad (7)$$

本文定义损失函数为

$$l(y, \mathbf{x}, \boldsymbol{\theta}) = \frac{1}{2}(y - f(\mathbf{x}, \boldsymbol{\theta}))^2. \quad (8)$$

对在线学习情形, 一次训练一个样本, 则在 t 时刻, 若训练样本为 (y_t, \mathbf{x}_t) , 参数 $\boldsymbol{\theta}_t$ 的迭代公式为

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \frac{\partial l(y_t, \mathbf{x}_t, \boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}_t}, \quad (9)$$

其中 η 是学习率.

$\frac{\partial l(y_t, \mathbf{x}_t, \boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}_t}$ 可表示为

$$\frac{\partial l(y_t, \mathbf{x}_t, \boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}_t} = \left\langle \frac{\partial l(y_t, \mathbf{x}_t, \boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}_t} \right\rangle + \delta_t, \quad (10)$$

其中: $\langle \cdot \rangle$ 表示期望, 相应的教师分布函数 $p_0(y, \mathbf{x})$ 为

$$p_0(y, \mathbf{x}) = q(\mathbf{x}) \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - f_0(\mathbf{x}))^2\right), \quad (11)$$

δ_t 为依赖于 $\boldsymbol{\theta}_t$, 均值为0的随机向量.

而对批处理学习情形, t 时刻参数 $\boldsymbol{\theta}_t$ 的迭代公式为

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \frac{1}{N} \sum_{i=1}^N \frac{\partial l(y_i, \mathbf{x}_i, \boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}_t}. \quad (12)$$

综合式(9)和式(12), 本文研究如下的平均方程:

$$\dot{\boldsymbol{\theta}}_{t+1} = \boldsymbol{\theta}_t - \eta \left\langle \frac{\partial l(y_t, \mathbf{x}_t, \boldsymbol{\theta}_t)}{\partial \boldsymbol{\theta}_t} \right\rangle. \quad (13)$$

当输入输出 (\mathbf{x}_t, y_t) 遍历历时, 式(12)可写为连续时间的形式, 即

$$\dot{\boldsymbol{\theta}} = -\eta \left\langle \frac{\partial l(y, \mathbf{x}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\rangle = -\eta \frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \quad (14)$$

其中 $L(\boldsymbol{\theta}) = \langle l(y, \mathbf{x}, \boldsymbol{\theta}) \rangle$.

显然, 当 N 比较大时, 式(12)可以近似等价于式(13). 所以式(14)不仅可以反映在线学习的特性, 也可以反映批处理学习的特性. 故在本文主要研究平均学

习方程(14). 平均损失函数 $L(\theta)$ 也是教师函数和学生函数之间的泛化误差.

3 平均学习方程的解析表达式(A analytical expression of averaged learning equations)

由于平均学习方程(14)既能反映在线学习的特性, 也能反映批处理学习的特性, 对研究MLPs的学习动态具有重要的作用, 所以给出平均学习方程的解析表达式具有重要的意义. 双曲正切函数 $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ 是一种双极性的激活函数, 在MLPs中具有广泛的应用. 然而其积分难以求解使得求取平均学习方程非常困难, 所以本文选取误差函数 $\phi(x) = \sqrt{\frac{2}{\pi}} \cdot \int_0^x \exp(-\frac{t^2}{2}) dt$ 作为MLPs的激活函数. 误差函数也是一种双极性的S型函数, 具有易积分的特点. 选取误差函数作为激活函数能够大大简化计算, 得到更丰富的理论分析成果. 基于此, 本文可得到平均学习方程的解析表达式.

定理 1 若选取误差函数作为MLPs的激活函数, 则平均学习方程(14)的解析表达式为

$$\dot{\mathbf{J}}_i = \eta w_i (\sum_{j=1}^2 v_j I_2(\mathbf{t}_j, \mathbf{J}_i) - \sum_{j=1}^2 w_j I_2(\mathbf{J}_j, \mathbf{J}_i)), \quad (15)$$

$$\dot{w}_i = \eta (\sum_{j=1}^2 v_j I_1(\mathbf{t}_j, \mathbf{J}_i) - \sum_{j=1}^2 w_j I_1(\mathbf{J}_j, \mathbf{J}_i)), \quad (16)$$

其中: $i = 1, 2$,

$$I_1(\mathbf{t}, \mathbf{J}) = \frac{2}{\pi} \arcsin \frac{\mathbf{t}^T \mathbf{J}}{\sqrt{1 + \|\mathbf{t}\|^2} \sqrt{1 + \|\mathbf{J}\|^2}}, \quad (17)$$

$$I_2(\mathbf{t}, \mathbf{J}) = \frac{2}{\pi} \sqrt{\det(\mathbf{B}(\mathbf{t}, \mathbf{J})^{-1})} \mathbf{A}^{-1} \mathbf{t}, \quad (18)$$

$$\mathbf{A} = \mathbf{I}_n + \mathbf{J} \mathbf{J}^T, \quad (19)$$

$$\mathbf{B}(\mathbf{t}, \mathbf{J}) = \mathbf{A} + \mathbf{t} \mathbf{t}^T = \mathbf{I}_n + \mathbf{J} \mathbf{J}^T + \mathbf{t} \mathbf{t}^T, \quad (20)$$

$$\mathbf{A}^{-1} = (\mathbf{I}_n + \mathbf{J} \mathbf{J}^T)^{-1} = \mathbf{I}_n - \frac{\mathbf{J} \mathbf{J}^T}{1 + \|\mathbf{J}\|^2}, \quad (21)$$

$$\mathbf{B}(\mathbf{t}, \mathbf{J})^{-1} = (\mathbf{I}_n - \frac{\mathbf{A}^{-1} \mathbf{t} \mathbf{t}^T}{1 + \mathbf{t}^T \mathbf{A}^{-1} \mathbf{t}}) \mathbf{A}^{-1}, \quad (22)$$

$$\det(\mathbf{B}(\mathbf{t}, \mathbf{J})^{-1}) = \frac{1}{(1 + \|\mathbf{t}\|^2)(1 + \|\mathbf{J}\|^2) - (\mathbf{t}^T \mathbf{J})^2}. \quad (23)$$

证 证明过程见附录A.

由定理1的结果, 本文还可以得出泛化误差 $L(\theta)$ 的解析表达式为

$$L(\theta) = \frac{1}{2} \sum_{i,j} v_i v_j I_1(\mathbf{t}_i, \mathbf{t}_j) - \sum_{i,j} v_i w_j I_1(\mathbf{t}_i, \mathbf{J}_j) + \frac{1}{2} \sum_{i,j} w_i w_j I_1(\mathbf{J}_i, \mathbf{J}_j). \quad (24)$$

4 互反奇异性区域附近的解析学习轨迹 (Analytical trajectories near opposite singularity)

4.1 坐标变换(Coordinate transformation)

为了方便进行理论分析, 本文首先引入与文献[15]研究重合奇异性相类似的坐标变换:

$$\mathbf{u} = \mathbf{J}_1 + \mathbf{J}_2, \quad (25)$$

$$\mathbf{v} = \frac{w_2 \mathbf{J}_2 + w_1 \mathbf{J}_1}{w_2 - w_1}, \quad (26)$$

$$w = w_2 - w_1, \quad (27)$$

$$z = \frac{w_2 + w_1}{w_2 - w_1}, \quad (28)$$

则原坐标可表示为

$$\mathbf{J}_1 = \frac{1}{2}(1+z)\mathbf{u} - \mathbf{v}, \quad (29)$$

$$\mathbf{J}_2 = \frac{1}{2}(1-z)\mathbf{u} + \mathbf{v}, \quad (30)$$

$$w_1 = \frac{1}{2}(z-1)w, \quad (31)$$

$$w_2 = \frac{1}{2}(z+1)w. \quad (32)$$

在新的坐标系统中模型参数为 $\xi = \{\mathbf{v}, w, \mathbf{u}, z\}$, 互反奇异性区域可表示为 $\mathcal{R} = \{\xi | \mathbf{u} = \mathbf{0}\}$. 显然 $z = 1$ ($z = -1$)表示 $w_1 = 0$ ($w_2 = 0$), 即文献[15]所定义的消去奇异性区域.

此时学生模型可表示为

$$f(\mathbf{x}, \xi) = \frac{1}{2}(z-1)w\phi(\mathbf{x}, \frac{1}{2}(1+z)\mathbf{u} - \mathbf{v}) + \frac{1}{2}(z+1)w\phi(\mathbf{x}, \frac{1}{2}(1-z)\mathbf{u} + \mathbf{v}). \quad (33)$$

对 $f(\mathbf{x}, \xi)$ 在 $\mathbf{u} = \mathbf{0}$ 处做泰勒展开:

$$f(\mathbf{x}, \xi) = w\phi(\mathbf{x}, \mathbf{v}) + \frac{1}{8}w(1-z^2)\mathbf{u}^T \frac{\partial^2 \phi(\mathbf{x}, \mathbf{v})}{\partial \mathbf{v} \partial \mathbf{v}^T} \mathbf{u} + O(\mathbf{u}^3). \quad (34)$$

此时学习方程(14)可重新表示为

$$\dot{\xi} = -\eta \mathbf{T} \mathbf{T}^T \left\langle \frac{\partial l(y, \mathbf{x}, \xi)}{\partial \xi} \right\rangle, \quad (35)$$

其中

$$\mathbf{T} = \frac{\partial \xi}{\partial \theta^T} = \begin{bmatrix} z-1 & z+1 & z+1 & 1-z \\ 2 & 2w & 2 & 2w \\ 0 & -1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & z+1 & 0 & 1-z \\ & w & & w \end{bmatrix} \mathbf{u} \quad (36)$$

为坐标变换的Jacobi矩阵.

由式(35)可得

$$\dot{\mathbf{v}} = \frac{z^2+1}{2} l_v + \frac{z^2+1}{2w^2} \mathbf{u} \mathbf{u}^T l_v -$$

$$\frac{z}{w}\mathbf{u}l_w + zl_u + \frac{z^2+1}{w^2}\mathbf{u}l_z, \quad (37)$$

$$\dot{w} = -\frac{z}{w}\mathbf{u}^T l_v + 2l_w - \frac{2z}{w}l_z, \quad (38)$$

$$\dot{\mathbf{u}} = zl_v + 2l_u, \quad (39)$$

$$\dot{z} = \frac{z^2+1}{w^2}\mathbf{u}^T l_v - \frac{2z}{w}l_w + \frac{2(z^2+1)}{w^2}l_z, \quad (40)$$

其中:

$$l_v = w\langle e(y, \mathbf{x}, \boldsymbol{\xi}) \frac{\partial \phi(\mathbf{x}, \mathbf{v})}{\partial \mathbf{v}} \rangle + O(\mathbf{u}^2), \quad (41)$$

$$l_w = \langle e(y, \mathbf{x}, \boldsymbol{\xi}) \phi(\mathbf{x}, \mathbf{v}) \rangle + \frac{1}{8}(1-z^2) \times \langle e(y, \mathbf{x}, \boldsymbol{\xi}) \mathbf{u}^T \frac{\partial^2 \phi(\mathbf{x}, \mathbf{v})}{\partial \mathbf{v} \partial \mathbf{v}^T} \mathbf{u} \rangle + O(\mathbf{u}^3), \quad (42)$$

$$l_u = \frac{1}{4}w(1-z^2)\langle e(y, \mathbf{x}, \boldsymbol{\xi}) \frac{\partial^2 \phi(\mathbf{x}, \mathbf{v})}{\partial \mathbf{v} \partial \mathbf{v}^T} \mathbf{u} \rangle + O(\mathbf{u}^2), \quad (43)$$

$$l_z = -\frac{1}{4}wz\langle e(y, \mathbf{x}, \boldsymbol{\xi}) \mathbf{u}^T \frac{\partial^2 \phi(\mathbf{x}, \mathbf{v})}{\partial \mathbf{v} \partial \mathbf{v}^T} \mathbf{u} \rangle + O(\mathbf{u}^3) \quad (44)$$

是 $L(\boldsymbol{\xi})$ 关于 $\boldsymbol{\xi}$ 的负梯度, $e(y, \mathbf{x}, \boldsymbol{\xi}) = f_0(\mathbf{x}) - f(\mathbf{x}, \boldsymbol{\xi})$ 是误差.

4.2 互反奇异性区域附近的渐近学习轨迹 (Asymptotical learning trajectories near opposite singularity)

当学习轨迹到达互反奇异性区域, 即 $\mathbf{u} = \mathbf{0}$ 附近时, 学习过程可看作为一个隐节点的学生模型去逼近两个隐节点的教师模型. 此时教师模型和学生模型可表示为如下的形式:

$$f(\mathbf{x}, \boldsymbol{\theta}_0) = v_1\phi(\mathbf{x}, \mathbf{t}_1) + v_2\phi(\mathbf{x}, \mathbf{t}_2), \quad (45)$$

$$f(\mathbf{x}, \boldsymbol{\theta}^*) = w^*\phi(\mathbf{x}, \mathbf{v}^*). \quad (46)$$

(\mathbf{v}^*, w^*) 表示教师模型的最佳逼近. 此时互反奇异性区域就可表示为

$$\mathcal{R}^* = \{\boldsymbol{\xi} | \mathbf{v} = \mathbf{v}^*, w = w^*, \mathbf{u} = \mathbf{0}, z \in \mathbb{R}\}. \quad (47)$$

在 \mathcal{R}^* 附近的学习动态可由 l_v, l_w, l_u 和 l_z 中 \mathbf{u} 的高阶项所得, 则式(41)–(44)变为

$$l_v(\boldsymbol{\xi}^*) = O(\mathbf{u}^2), \quad (48)$$

$$l_w(\boldsymbol{\xi}^*) = \frac{1}{2} \frac{1-z^2}{w^{*2}} \mathbf{u}^T H(\mathbf{v}^*, w^*) \mathbf{u} + O(\mathbf{u}^3), \quad (49)$$

$$l_u(\boldsymbol{\xi}^*) = (1-z^2)H(\mathbf{v}^*, w^*) \mathbf{u} + O(\mathbf{u}^2), \quad (50)$$

$$l_z(\boldsymbol{\xi}^*) = -z\mathbf{u}^T H(\mathbf{v}^*, w^*) \mathbf{u} + O(\mathbf{u}^3), \quad (51)$$

其中

$$H(\mathbf{v}^*, w^*) = \frac{1}{4}w^*\langle e(y, \mathbf{x}, \boldsymbol{\xi}) \frac{\partial^2 \phi(\mathbf{x}, \mathbf{v})}{\partial \mathbf{v} \partial \mathbf{v}^T} \rangle |_{\boldsymbol{\xi}=\boldsymbol{\xi}^*}. \quad (52)$$

此时 $l_u(\boldsymbol{\xi}^*)$ 是 \mathbf{u} 的同阶无穷小量, $l_v(\boldsymbol{\xi}^*), l_w(\boldsymbol{\xi}^*)$ 和 $l_z(\boldsymbol{\xi}^*)$ 都是 \mathbf{u}^2 的同阶无穷小量. 本文只关注 \mathbf{u} 和 z 的变化, 则对 \mathcal{R}^* 附近的学习动态, 式(39)和式(40)可重新

表示为

$$\dot{\mathbf{u}} = 2(1-z^2)H(\mathbf{v}^*, w^*)\mathbf{u}, \quad (53)$$

$$\dot{z} = -\frac{z(1-z^2)}{w^{*2}}\mathbf{u}^T H(\mathbf{v}^*, w^*)\mathbf{u} - \frac{2z(z^2+1)}{w^{*2}}\mathbf{u}^T H(\mathbf{v}^*, w^*)\mathbf{u}. \quad (54)$$

定理 2 由式(53)和式(54), 可得到 \mathcal{R}^* 附近的平均学习方程的轨迹为

$$h(\mathbf{u}) = \frac{1}{2}\mathbf{u}^T \mathbf{u} = \frac{2w^{*2}}{3} \log \frac{(z^2+3)^2}{|z|} + C, \quad (55)$$

其中 $h(\mathbf{u}) = \frac{1}{2}\mathbf{u}^T \mathbf{u}$, C 为由初始模型参数 $(u^{(0)}, z^{(0)})$ 决定的常数. 轨迹如图2所示.

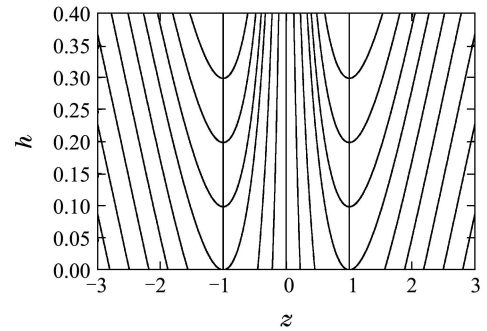


图2 \mathcal{R}^* 附近的理论学习轨迹

Fig. 2 Theoretical learning trajectories near \mathcal{R}^*

证 令 $h(\mathbf{u}) = \frac{1}{2}\mathbf{u}^T \mathbf{u}$, 由式(53)和式(54)可得

$$\dot{h} = \mathbf{u}^T \dot{\mathbf{u}} = \frac{2w^{*2}(z^2-1)}{z(z^2+3)}\dot{z}. \quad (56)$$

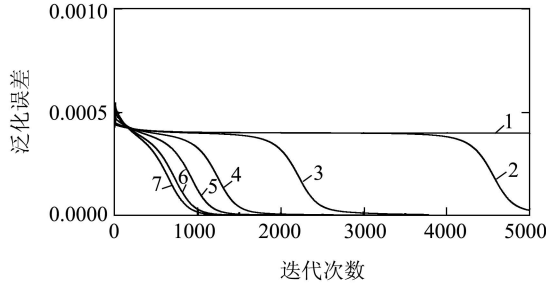
通过解微分方程(56), 可得到 \mathcal{R}^* 附近的平均学习方程的轨迹为式(55)中的形式.

4.3 实际平均学习轨迹(Real averaged learning trajectories)

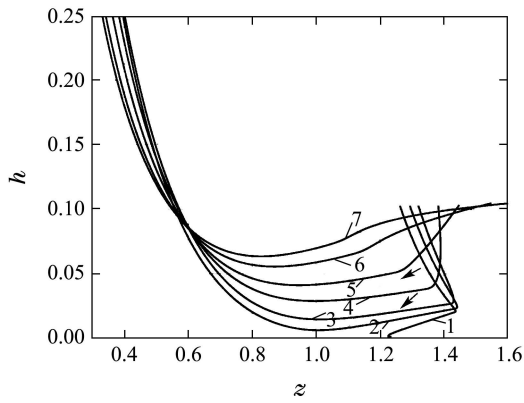
在给出互反奇异性区域附近的理论学习轨迹后, 笔者现在研究MLPs在学习过程中互反奇异性区域附近实际的学习轨迹. 由于平均学习方程(14)是常微分方程, 所以给定教师函数的权值和学生函数的初始权值后, 通过解微分方程就可以得出学生函数各个参数的学习轨迹, 进一步就可以得到 h 和 z 的学习轨迹. 本文取教师函数的权值为 $\mathbf{t}_1 = (0.2, -0.5)^T$, $w_1 = 0.5$, $\mathbf{t}_2 = (0.7, 0.4)^T$, $w_2 = 0.3$. 给定学生函数的初始权值, 可以得出 $h \sim z$ 的轨迹. 图3表示由不同的学生函数初始权值得出的泛化误差轨迹和相应的 $h \sim z$ 轨迹.

对比图2和图3(b), 可以看出图3(b)和图2的右半部分比较相似, 而随着 h 越来越大时, 两图之间的差别越来越大. 这是因为理论学习轨迹是通过在 $\mathbf{u} = \mathbf{0}$, 即 $h = 0$ 处做泰勒展开求得, 所以在 h 比较小时能保证足够的精度, 而当 h 越来越大时, 误差也会越来越大. 由

图3(b)的曲线1可以看出在训练之后 h 接近为0,即两个隐节点接近互反,此时由图3(a)中的曲线1可以看出泛化误差基本不变.图3(b)中的曲线2-7表明学习轨迹穿越 $z = 1$ 即消去奇异性区域,对应图3(a)中的曲线可以看出误差曲线发生了明显的平坦区现象,最终到达最佳逼近.



(a) 泛化误差轨迹



(b) $h \sim z$ 轨迹

图3 \mathcal{R}^* 附近的实际学习轨迹

Fig. 3 Real learning trajectories near \mathcal{R}^*

5 仿真例子(Simulation examples)

本节对平均学习动态(14)、批处理学习动态(12)和在线学习动态(9)进行比较.通过解微分方程(15)和式(16),可以用数值方法求得MLPs的平均学习动态.批处理学习动态和在线学习动态通过仿真实验得出.选取如下的教师函数:

$$f_0(\mathbf{x}) = v_1\phi(\mathbf{x}, \mathbf{t}_1) + v_2\phi(\mathbf{x}, \mathbf{t}_2), \quad (57)$$

其中: $w_1 = 0.5, \mathbf{t}_1 = (0.2, -0.5)^T, w_2 = 0.3, \mathbf{t}_2 = (0.7, 0.4)^T$.

5.1 平均学习动态(Averaged learning dynamics)

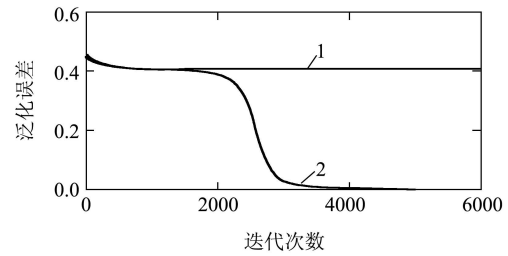
本文首先研究MLPs的平均学习动态.除了泛化误差这个重要的指标外,由于 h 接近0时可表示两个隐节点接近互反,所以 h 的轨迹也具有重要的意义.表1和表2为学生函数的初始权值和通过解平均学习方程得到的最终权值.

图4中的曲线1和曲线2表示当学生函数初始权值分别取表1和表2中的值时泛化误差 $L(\theta)$ 和 h 的学习

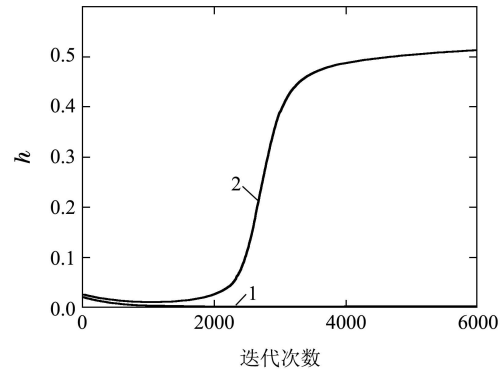
轨迹.通过图4(b)中的曲线1,可以看到 h 接近于0,即学生函数的两个隐节点接近互反,此时学生参数进入互反奇异性区域.由图4(a)中的曲线1可以看出泛化误差基本不变,学习过程受到互反奇异性区域的严重影响.由表1的最终权值可以看出两个隐节点并没有完全互反,但此时

$$\frac{\partial L(\theta)}{\partial \theta} = 10^{-4} \times (0.1655, 0.0095, -0.1066, -0.1057, 0.0501, 0.1223)^T$$

已经很小,说明即使经过更长的运算,两个隐节点已经基本停止迭代,学生参数也很难离开互反奇异性区域,学生函数难以达到全局最优.



(a) 泛化误差轨迹



(b) h 的轨迹

图4 平均学习动态

Fig. 4 Averaged learning dynamics

表2中的初始权值仅是对表1中的初始权值做一个小的改动,由图4(a)中的曲线2可以看出 $L(\theta)$ 有一段时间基本不变,然后突然快速下降,发生了平坦区现象.由表2的最终权值可以看出虽然 $\mathbf{J}_2 = (-0.1811, 0.4943)^T, w_2 = -0.5179$ 与教师参数 $\mathbf{t}_1 = (0.2, -0.5)^T, w_1 = 0.5$ 符号不同,但由于激活函数为双极性函数, $w_2\phi(\mathbf{x}, \mathbf{J}_2) = -v_1\phi(\mathbf{x}, -\mathbf{t}_1) = v_1\phi(\mathbf{x}, \mathbf{t}_1)$,即学生函数已经到达了全局最优点,实现了最佳逼近.

表1 学习过程陷入互反奇异性区域

Table 1 The learning process trapped in the opposite singularities

	初始 \mathbf{J}	初始 w	最终 \mathbf{J}	最终 w
隐节点1	(0.40, 0.10)	-0.12	(0.2901, -0.1457)	-0.1152
隐节点2	(-0.50, 0.35)	-1.08	(-0.3309, 0.1667)	-0.9129

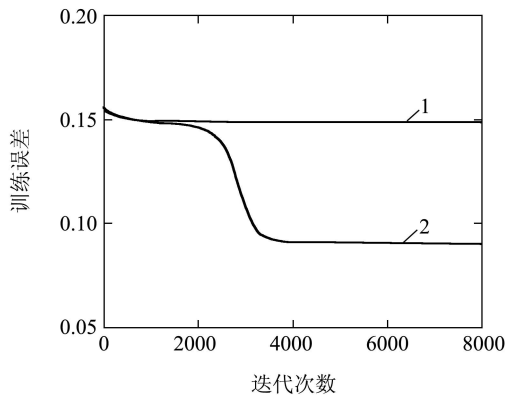
表 2 学习过程受奇异性影响发生平坦区现象

Table 2 Plateau occurred in the learning process affected by singularities

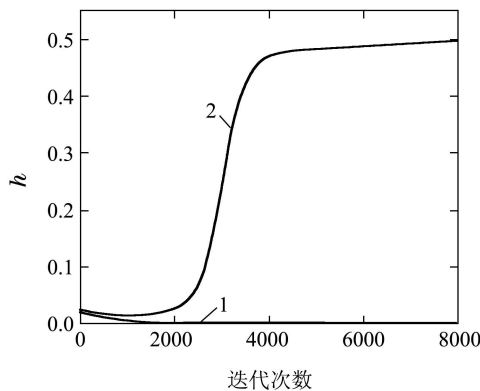
	初始 J	初始 w	最终 J	最终 w
隐节点1	(0.40, 0.10)	-0.12	(0.6672, 0.3947)	0.3199
隐节点2	(-0.50, 0.35)	-0.90	(-0.1811, 0.4943)	-0.5179

5.2 批处理学习动态(Batch mode learning dynamics)

下面本文研究批处理学习中MLPs的学习动态. 由教师函数(43)生成200个样本, 样本输入服从标准高斯分布, 噪声 ε 服从均值为0和方差为0.05的高斯分布. 本文使用梯度下降法进行训练, 学习率 $\eta = 0.005$, 迭代次数为8000. 在批处理学习中本文用训练误差代替泛化误差. 图5(a)和图5(b)分别表示训练误差和 h 的轨迹, 其中曲线1为学生参数取表1中的初始值, 曲线2为学生参数取表2中的初始值.



(a) 训练误差轨迹



(b) h 的轨迹

图 5 批处理学习动态

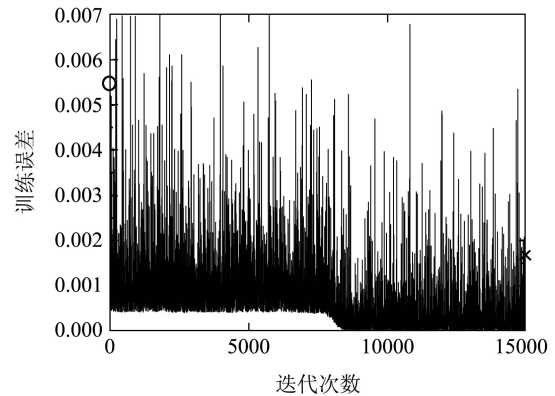
Fig. 5 Batch mode dynamics of learning

由于样本数和噪声的影响, 在批处理学习中训练误差达不到0. 而平均学习方程在计算时对教师分布函数取期望, 所以取期望后系统的平均学习动态就避免了样本数和噪声的影响. 对比平均学习动态和批处

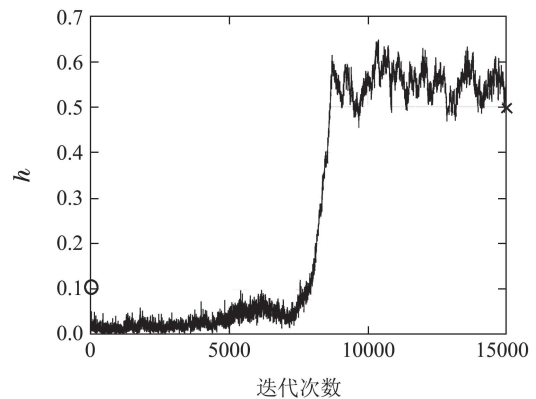
理学习动态的泛化误差轨迹和 h 的学习轨迹, 可以发现两者的轨迹十分相似, 而且训练结束后由平均学习方程训练得到的模型最终权值和由批处理模式训练得到的模型最终权值基本相同, 因此本文可以使用平均学习方程来研究MLPs的批处理学习动态. 使用平均学习方程可避免样本和噪声的影响, 并且学习速度更快, 训练结果也更稳定.

5.3 在线学习动态(Online learning dynamics)

现在本文研究在线学习情形下MLPs的学习动态. 批处理学习模式为所有样本同时训练, 而在线学习模式为一次训练一个样本. 样本的生成方式与批处理学习模式中相同. 笔者生成15000个样本, 每个样本训练200次, 学习率 $\eta = 0.05$. 图6(a)和图6(b)分别表示训练误差和 h 的轨迹, 学生参数初始值取表1中的初始值. “o”和“x”分别表示在线学习过程中的初始状态和最终状态.



(a) 训练误差轨迹



(b) h 的轨迹

图 6 在线学习动态

Fig. 6 Online dynamics of learning

由图6的曲线可以看出由于样本的随机性和噪声的存在, 训练过程一直波动. 对比图4和图6的曲线可以看出在线学习过程中的学习曲线是沿着平均学习过程的学习曲线波动. 由图6(b)可以看出在训练开始后 h 迅速下降, 接近于0, 学习过程在重合奇异性区域附近波动, 学习过程受到奇异性区域的影响. 由图

4和图5的仿真结果可以看出平均学习模式和批处理学习模式中学习过程受到重合奇异性区域影响后会最终陷入重合奇异性区域, 学生模型到达局部最优值. 这是由于平均学习动态计算过程中经过对教师函数取期望后不再受样本输入随机性和附加噪声的影响, 而批处理模式中一次学习多个样本, 也会减少这种随机性的影响, 所以平均学习动态和批处理学习动态会陷入互反奇异性区域. 而在线学习模式的情形不同, 由于一次只训练一个样本, 样本的随机性和附加噪声的影响使得学习过程在重合奇异性区域附近随机游走, 由图6可以看出学习过程最终离开重合奇异性区域从而达到最优值, 发生了明显的平坦区现象.

6 结论(Conclusions)

多层感知器神经网络(MLPs)是神经网络中应用最广泛的前馈神经网络之一, 然而在用BP算法训练时经常出现学习速度很慢, 容易陷入极小点等奇异行为. 经研究发现, 这些奇异行为都是由参数空间中存在的奇异性区域导致的. 在这些奇异性区域, Fisher信息阵退化, 导致学习困难, 学习过程变得很慢. 当MLPs的两个隐节点权值接近互反时, Fisher信息阵接近奇异, 学习过程受到很大影响, 因此很有必要研究互反奇异性区域附近的学习动态. 为了克服理论分析中传统激活函数 $\tan(x)$ 积分困难, 本文选取误差函数作为MLPs的激活函数, 并且得到了MLPs的平均学习方程的解析表达式. 引入坐标变换和进行泰勒展开, 得到了MLPs在互反奇异性附近的理论学习轨迹, 并通过使用平均学习方程得到了实际的学习轨迹, 进行了对比分析. 在仿真实验中本文分别对MLPs的平均学习动态、批处理学习动态和在线学习动态进行了比较分析. 仿真结果表明可以使用平均学习方程来研究MLPs的批处理学习动态.

参考文献(References):

- [1] WATANABE S. Almost all learning machines are singular [C] // *Proceedings of IEEE Symposium on Foundations of Computational Intelligence*. Honolulu, HI: IEEE, 2007: 383 – 388.
- [2] FUKUMIZU K. A regularity condition of the information matrix of a multilayer perceptron network [J]. *Neural Networks*, 1996, 9(5): 871 – 879.
- [3] FUKUMIZU K, AMARI S. Local minima and plateaus in hierarchical structure of multilayer perceptrons [J]. *Neural Networks*, 2000, 13(3): 317 – 327.
- [4] AMARI S, PARK H, OZEKI Z. Singularities affect dynamics of learning in neuromanifolds [J]. *Neural Computation*, 2006, 18(5): 1007 – 1065.
- [5] SAAD D, SOLLA A. Exact solution for online learning in multilayer neural networks [J]. *Physical Review Letters*, 1995, 74(21): 4337 – 4340.
- [6] FREEMAN J A S, SAAD D. Dynamics of on-line learning in radial basis function networks [J]. *Physical Review E*, 1997, 56(1): 907 – 918.

- [7] AMARI S, NAGAOKA H. *Information Geometry* [M]. New York: AMS and Oxford University Press, 2000.
- [8] PARK H, INOUE M, OKADA M. Online learning dynamics of multilayer perceptrons with unidentifiable parameters [J]. *Journal of Physics A: Mathematical and General*, 2003, 36(47): 11753 – 11764.
- [9] AMARI S, OZEKI T, COUSSEAU F, et al. Dynamics of learning in hierarchical models—singularity and milnor attractor [C] // *Advances in Cognitive Neurodynamics (II)*. Hangzhou: Springer Netherlands, 2011: 3 – 9.
- [10] 魏海坤, 李奇, 宋文忠. 梯度算法下RBF网的参数变化动态 [J]. *控制理论与应用*, 2007, 24(3): 356 – 360.
(WEI Haikun, LI Qi, SONG Wenzhong. Gradient learning dynamics of radial basis function networks [J]. *Control Theory & Applications*, 2007, 24(3): 356 – 360.)
- [11] WATANABE S. Algebraic geometrical methods for hierarchical learning machines [J]. *Neural Networks*, 2001, 14(8): 1049 – 1060.
- [12] WATANABE S. A widely applicable bayesian information criterion [J]. *Journal of Machine Learning Research*, 2013, 14(3): 867 – 897.
- [13] PARK H, OZEKI T. Singularity and slow convergence of the EM algorithm for gaussian mixtures [J]. *Neural Process Letters*, 2009, 29(1): 45 – 59.
- [14] FREEMAN J A S, SAAD D. On-line learning in radial basis function networks [J]. *Neural Computation*, 1997, 9(7): 1601 – 1622.
- [15] WEI H, ZHANG J, COUSSEAU F, et al. Dynamics of learning near singularities in layered networks [J]. *Neural Computation*, 2008, 20(3): 813 – 843.
- [16] WEI H, AMARI S. Dynamics of learning near singularities in radial basis function networks [J]. *Neural Networks*, 2008, 21(7): 989 – 1005.
- [17] COUSSEAU F, OZEKI T, AMARI S. Dynamics of learning in multi-layer perceptrons near singularities [J]. *IEEE Transactions on Neural Networks*, 2008, 19(8): 1313 – 1328.

附录 A(Appendix A)

由式(6)可得

$$y - f_0(\mathbf{x}) = \varepsilon \sim \mathcal{N}(0, 1), \quad (\text{A1})$$

则有

$$\begin{aligned} & \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} \exp\left(-\frac{1}{2}(y - f_0(\mathbf{x}))^2\right) dy = \\ & \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} \exp\left(-\frac{\varepsilon^2}{2}\right) d\varepsilon = 1. \end{aligned} \quad (\text{A2})$$

$I_1(\mathbf{s}, \mathbf{v})$ 和 $I_2(\mathbf{s}, \mathbf{v})$ 可重新表示为

$$\begin{aligned} I_1(\mathbf{s}, \mathbf{v}) &= (2\pi)^{-\frac{n}{2}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \phi(\mathbf{s}^T \mathbf{x}) \phi(\mathbf{v}^T \mathbf{x}) e^{(-\frac{1}{2}\|\mathbf{x}\|^2)} \cdot \\ & \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - f_0(\mathbf{x}))^2\right) dy d\mathbf{x} = \\ & (2\pi)^{-\frac{n}{2}} \int_{-\infty}^{\infty} \phi(\mathbf{s}^T \mathbf{x}) \phi(\mathbf{v}^T \mathbf{x}) e^{(-\frac{1}{2}\|\mathbf{x}\|^2)} d\mathbf{x}, \end{aligned} \quad (\text{A3})$$

$$\begin{aligned} I_2(\mathbf{s}, \mathbf{v}) &= (2\pi)^{-\frac{n}{2}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \phi(\mathbf{s}^T \mathbf{x}) \frac{\partial \phi(\mathbf{v}^T \mathbf{x})}{\partial \mathbf{v}} e^{(-\frac{1}{2}\|\mathbf{x}\|^2)} \cdot \\ & \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - f_0(\mathbf{x}))^2\right) dy d\mathbf{x} = \\ & (2\pi)^{-\frac{n}{2}} \int_{-\infty}^{\infty} \phi(\mathbf{s}^T \mathbf{x}) \frac{\partial \phi(\mathbf{v}^T \mathbf{x})}{\partial \mathbf{v}} e^{(-\frac{1}{2}\|\mathbf{x}\|^2)} d\mathbf{x}. \end{aligned} \quad (\text{A4})$$

由莱布尼兹积分法则可得

$$I_2(\mathbf{s}, \mathbf{v}) = \frac{\partial I_1(\mathbf{s}, \mathbf{v})}{\partial \mathbf{v}}, \quad (\text{A5})$$

$$\begin{aligned}
I_2(\mathbf{s}, \mathbf{v}) &= \\
(2\pi)^{-\frac{n}{2}} \int_{-\infty}^{\infty} \phi(\mathbf{s}^T \mathbf{x}) \frac{\partial \phi(\mathbf{v}^T \mathbf{x})}{\partial \mathbf{v}} e^{(-\frac{1}{2} \|\mathbf{x}\|^2)} d\mathbf{x} &= \\
(2\pi)^{-\frac{n}{2}} \frac{2}{\pi} \int_{-\infty}^{\infty} \phi(\mathbf{s}^T \mathbf{x}) \mathbf{x} e^{(-\frac{1}{2} (\mathbf{v}^T \mathbf{x})^2)} e^{(-\frac{1}{2} \|\mathbf{x}\|^2)} d\mathbf{x} &= \\
(2\pi)^{-\frac{n}{2}} \frac{2}{\pi} \int_{-\infty}^{\infty} \mathbf{x} \phi(\mathbf{s}^T \mathbf{x}) e^{(-\frac{1}{2} (\|\mathbf{x}\|^2 + (\mathbf{v}^T \mathbf{x})^2))} d\mathbf{x} &= \\
\frac{2}{\pi} \sqrt{\det(\mathbf{B}^{-1})} \mathbf{A}^{-1} \mathbf{s}, & \quad (\text{A6})
\end{aligned}$$

其中:

$$\mathbf{A} = \mathbf{I}_n + \mathbf{v}\mathbf{v}^T, \quad (\text{A7})$$

$$\mathbf{B} = \mathbf{A} + \mathbf{s}\mathbf{s}^T. \quad (\text{A8})$$

根据Sherman-Morrison公式, 可得

$$\mathbf{A}^{-1} = (\mathbf{I}_n + \mathbf{v}\mathbf{v}^T)^{-1} = \mathbf{I}_n - \frac{\mathbf{v}\mathbf{v}^T}{1 + \|\mathbf{v}\|^2}, \quad (\text{A9})$$

$$\mathbf{B}^{-1} = (\mathbf{A} + \mathbf{s}\mathbf{s}^T)^{-1} = (\mathbf{I}_n - \frac{\mathbf{A}^{-1}\mathbf{s}\mathbf{s}^T}{1 + \mathbf{s}^T \mathbf{A}^{-1} \mathbf{s}}) \mathbf{A}^{-1}. \quad (\text{A10})$$

使用Sylvester行列式定理可得

$$\det(\mathbf{A}^{-1}) = 1 - \frac{\|\mathbf{v}\|^2}{1 + \|\mathbf{v}\|^2} = \frac{1}{1 + \|\mathbf{v}\|^2}, \quad (\text{A11})$$

$$\det(\mathbf{B}^{-1}) = \frac{1}{(1 + \|\mathbf{s}\|^2)(1 + \|\mathbf{v}\|^2) - (\mathbf{s}^T \mathbf{v})^2}. \quad (\text{A12})$$

由式(A5), I_1 可通过 I_2 关于 \mathbf{v} 积分求得

$$\begin{aligned}
I_1(\mathbf{s}, \mathbf{v}) &= \int_{-\infty}^{\mathbf{v}} I_2(\mathbf{s}, \mathbf{v}) d\mathbf{v} = \\
\frac{2}{\pi} \int_{-\infty}^{\mathbf{v}} \frac{\mathbf{A}^{-1} \mathbf{s}}{\sqrt{(1 + \|\mathbf{s}\|^2)(1 + \|\mathbf{v}\|^2) - (\mathbf{s}^T \mathbf{v})^2}} d\mathbf{v} &= \\
\frac{2}{\pi} (\arcsin \frac{\mathbf{s}^T \mathbf{v}}{\sqrt{1 + \|\mathbf{s}\|^2} \sqrt{1 + \|\mathbf{v}\|^2}} + C), & \quad (\text{A13})
\end{aligned}$$

其中 C 是一个常数.

由式(A3)可知 $I_1(\mathbf{0}, \mathbf{0}) = 0$, 故 $C = 0$. 由此可得

$$I_1(\mathbf{s}, \mathbf{v}) = \frac{2}{\pi} \arcsin \frac{\mathbf{s}^T \mathbf{v}}{\sqrt{1 + \|\mathbf{s}\|^2} \sqrt{1 + \|\mathbf{v}\|^2}}. \quad (\text{A14})$$

作者简介:

郭伟立 (1987-), 男, 博士研究生, 研究方向为前馈神经网络奇异学习动态, E-mail: weiliguo@seu.edu.cn;

魏海坤 (1971-), 男, 博士, 教授, 研究方向为流程工业综合自动化系统设计、实时系统设计、神经网络学习理论、模式识别和人工智能等, E-mail: hkwei@seu.edu.cn;

赵军圣 (1980-), 男, 博士研究生, 讲师, 研究方向为非线性系统奇异学习动态, E-mail: zhaojunshao@163.com;

张侃健 (1972-), 男, 博士, 教授, 研究方向为非线性系统的鲁棒控制、优化控制等, E-mail: kjzhang@seu.edu.cn.