

对偶线性规划神经网络及算法步长的选取

裴炳南, 保 铮

(西安电子科技大学雷达信号处理国家重点实验室, 陕西西安 710071)

摘 要: 文章在对偶线性规划框架内研究神经网络的分析性质, 用特征值方法界定了 Hess 阵和神经网络学习步长的取值范围, 给出了用李普希兹常数表示的算法步长公式. 数值仿真结果表明, 理论分析是正确的, 给出的算法步长公式是有效的.

关键词: 对偶线性规划; 人工神经网络; 最优化; 学习率; 计算机仿真

中图分类号: TP183 **文献标识码:** A **文章编号:** 0372-2112 (2002) 01-0110-04

Study of Dual Linear Programming Neural Networks and Selection of Its Learning Rate

PEI Bing-nan, BAO Zheng

(Key Lab. of Radar Signal Processing, Xidian University, Xi'an, Shaanxi 710071, China)

Abstract: The research into the properties of neural networks is made in the framework of the dual linear programming theory. The variation range of a Hessian matrix as well as that of learning rate of the networks is confined by means of the eigenvalue. A Lipschitz constant based formula for the algorithm's convergence is given. Simulation is given to illustrate that the theory is correct and the formula given is efficacious.

Key words: dual linear programming; artificial neural networks; optimization; learning rate; simulation

1 引言

线性规划方法在信号处理、机器人技术、自动控制、系统理论和统计学中有广泛应用^[1]. 解决线性规划的经典方法单纯形法, 该方法在大规模约束条件时计算量很大^[2], 在一些实时处理场合不敷应用. 另一种方法是人工神经网络方法. 由于该方法具有并行处理的特点, 人们渴望它能弥补单纯形法的不足, 并做了大量理论研究和仿真实验. 然而理论研究一直忽视或回避神经网络状态更新增益因子或算法步长的选取问题. 这个问题又是工程实践和计算机仿真中必须解决的问题; 因为实现神经网络的差分方程(离散时间域)存在算法步长的上界, 超过这个上界, 算法就不收敛.

在文献中, 除了纯数学的叙述和证明之外, 计算机仿真是说明神经网络的有效手段. 人们往往凭直觉或经验用试错法选取步长. 这样选取步长缺乏理论指导, 选定的步长没有多大实际意义, 用于小规模约束条件尚可, 对于大规模约束问题的求解, 则难以凑效. 文献[3]将能量函数视为二次型, 利用泰勒级数展开研究了算法步长的选取问题. 可惜能量函数本身不是二次型的, 因而, 结果有一定局限性. 本文旨在从理论上给出基于对偶线性规划神经网络算法步长的理论界限和实用公式.

文章安排如下. 第二节提出基于对偶理论的线性规划神经网络模型, 第三节分析该模型的收敛性质, 第四节推导出算法实现的步长公式, 第五节给出了数值仿真结果, 最后是结论.

2 基于对偶理论的线性规划神经网络模型

2.1 理论根据

设线性规划问题的原规划是

$$(LP) \quad \begin{aligned} \max \quad & c^T x \\ \text{s. t.} \quad & Ax \leq b \\ & x \geq 0 \end{aligned} \quad (1)$$

其中, $A \in R^m \times n$, $b \in R^m$, $x \in R^n$, $c \in R^n$. 可行域 $R_p = \{x | Ax \leq b, x \geq 0, x \in R^n\}$.

原规划的对偶规划是

$$(DP) \quad \begin{aligned} \min \quad & b^T y \\ \text{s. t.} \quad & A^T y \leq c \\ & y \geq 0 \end{aligned} \quad (2)$$

其可行域是 $R_D = \{y | A^T y \leq c, y \geq 0, y \in R^m\}$.

下述线性规划对偶理论表明, 可利用原规划与对偶规划之间的关系, 建立线性规划神经网络的能量函数.

定理 1^[2] 若 $x^* \in R_p, y^* \in R_D$, 且 $c^T x^* = b^T y^*$, 则 x^* , y^* 分别是 (LP) 与 (DP) 的最优解.

定理 2^[2] 若同时存在容许解和对偶容许解, 则必同时存在基本最优解和基本对偶最优解, 且最优值相等, 即 $\max_{x \in R_p}$

$$c^T x = \min_{y \in R_D} b^T y.$$

2.2 神经网络能量函数的构造

记 $z = \begin{bmatrix} x \\ y \end{bmatrix}, f = \begin{bmatrix} -c \\ b \end{bmatrix}, l = \begin{bmatrix} 0 \\ -A \\ 0 \end{bmatrix}, J$, 合并 (1) (2) 两式写成

$$\begin{aligned} \min \quad & z^T z \\ \text{s.t.} \quad & z + f \geq 0 \\ & z \geq 0 \end{aligned} \quad (3)$$

为叙述方便, 引入下面符号定义. 设向量 $x = [x_1, x_2, x_3, \dots, x_n]^T$.

定义 1 称向量 $x > 0$, 如果 $x_1 > 0, x_2 > 0, \dots, x_n > 0$ 同时成立.

定义 2 $|x| = [|x_1|, |x_2|, \dots, |x_n|]^T$

定义 3 $x^- = \min(0, x) = \frac{1}{2}(x - |x|)$

定义 4 $x^+ = x^T x$.

定理 3 函数 x^- 是连续的.

证略.

定理 4 $x \geq 0$ 与 $x^- = 0$ 互等价.

证略.

利用上述概念, 可以构造对应于线性规划 (3) 的能量函数如下:

$$E(z) = \frac{1}{2} z^T z + \frac{1}{2} (z + f)^-{}^2 + \frac{1}{2} z^-{}^2 \quad (4)$$

显然 $E(z) \geq 0$, 等号当且仅当式 (4) 右边三项同时为零时成立. 由定理 (4) 知, $E(z) = 0$ 的解就是规划问题 (3) 的最优解, 换言之

$$\{z = (x^T, y^T)^T | E(z) = 0\} = \{z = (x^T, y^T)^T | z^T z = 0, z + f \geq 0, z \geq 0\}$$

方程 (4) 把对偶线性规划问题 (3) 映射为能量函数最小点的求解问题.

2.3 神经网络时间动力学方程建立

用经典单纯形法求解线性规划问题, 往往需要构造人工变量, 然后通过不断地换基运算求出最优解^[2]. 在用神经网络求解问题时, 状态 z 被理解为状态空间中一点, 如果约束方程是渐近稳定的, 即么, 随着时间的推移, 状态 z 将收敛到问题的解. $z = z(t)$ 是时间的函数.

考察能量函数 $E(z(t))$ 随时间演化的规律,

$$\frac{dE(z)}{dt} = \left(\frac{dz}{dt}\right)^T \nabla E(z) \quad (5)$$

采用广义梯度法^[1]取

$$\frac{dz}{dt} = -\mu(t) \nabla E(z) \quad (6)$$

为线性规划神经网络方程

仿附录引理 1 的求导方法, 不难得到

$$\nabla E(z) = z + (z + f)^- + z^- \quad (7)$$

3 神经网络的收敛性能

定理 5 能量函数 $E(z)$ 是可微不增凸函数. (证明见附录)

定理 6 梯度 $\nabla E(z)$ 是李普希兹 (Lipschitz) 连续的.

证 由附录引理 2 和方程 (7) 得

$$\begin{aligned} \nabla E(z_1) - \nabla E(z_2) &= z_1 + (z_1 + f)^- + z_1^- - [z_2 + (z_2 + f)^- + z_2^-] \\ &= (z_1 - z_2) + [(z_1 + f)^- - (z_2 + f)^-] + (z_1^- - z_2^-) \\ &= L \cdot (z_1 - z_2) \end{aligned}$$

这里 L 表示矩阵的某种范数运算. 对于给定的线性规划问题, 李普希兹常数

$$L = \|z + (z + f)^- + z^-\| \quad (8)$$

是一个确定的正数. 证毕

定理 7 设 $S_1 = \{z = (x^T, y^T)^T | R^{m+n} | \nabla E(z) = 0\}$ 是方程 (6) 的平稳点集, $S_2 = \{z^* = (x_0^T, y_0^T)^T | R^{m+n} | x_0, y_0\}$ 分别是原规划 (LP) 和对偶规划 (DP) 的最优解/是优化问题 (3) 的最优解集, 则这两个集合相等, $S_1 = S_2$.

证 (略)

定理 8 设原规划 (1) 与对偶规划 (2) 有唯一最优解 $z^* = [x_0^T, y_0^T]^T$, 则神经网络 (6) 全局一致渐近稳定收敛到该最优解. 当原规划与对偶规划有无穷多个最优解时, 则每个解都是李雅普诺夫稳定的.

证 (略)

4 离散时间循环神经网络及其收敛步长的选取

由微分方程式 (6) 描述的是循环 (Recurrent) 神经网络. 用数字电路实现将产生算法步长的选取问题. 这是神经网络工程实现时不能回避的问题, 也是计算机仿真时必须考虑的问题.

4.1 离散时间神经网络方程的导出

将式 (7) 代入式 (6), 将神经网络方程 (6) 写成

$$\frac{dz}{dt} = -\mu(t) [z + (z + f)^- + z^-] \quad (9)$$

$$z(0) = [x^T(0), y^T(0)]^T$$

假设微分方程 (6) 已经解出, 解的轨线为 $z(t)$. 在时间轴上均匀取点 $t_0 (= 0), t_1, t_2, \dots$, 由式 (9) 知解轨线至少是一阶光滑的, 因此, 在邻域内, $z(t_{j+1})$ 可用一阶泰勒级数展开为

$$z(t_{j+1}) = z(t_j) - \mu(t_j) t_j \nabla E(z(t_j)) \quad (10)$$

实用上取 $\mu(t_j)$ 为小的正数, 令 $\mu = \mu(t_j) t_j$, 并用符号 z_j 表示 $z(t_j)$, 重写 (10) 式为

$$z_{j+1} = z_j - \mu \nabla E(z_j) = z_j - \mu [z_j + (z_j + f)^- + z_j^-] \quad (11)$$

显然, 式 (11) 是典型的最速下降法公式. 众所周知, 在最速下降法中, 算法步长 μ 不仅控制算法收敛速率, 而且还关系着算法是否收敛. 理论上存在一个上界 μ^* , 当 $0 < \mu < \mu^*$, 算法 (11) 均方收敛. 下面研究 μ^* 的估计问题.

4.2 算法步长上界 μ^* 的估计及 μ 的确定



为了估计算法步长,我们先研究一下能量函数 $E(z)$ 的 Hess 矩阵 $\nabla^2 E(z)$ 的计算和估计问题.

应当指出,在经典微分学意义下, Hessian 阵 $H(z) = \nabla^2 E(z)$ 不存在. 为了分析神经网络 (11) 的收敛条件,我们用广义函数论中关于冲激函数的概念求 $H(z)$.

定义 5 设 R 是实矩阵, $R = (r_{ij})$, 则符号矩阵函数定义为

$$\text{sgn}(R) \triangleq (\text{Sgn}(r_{ij}))$$

其中,
$$\text{sgn}(r_{ij}) = \begin{cases} 1, & r_{ij} \geq 0 \\ -1, & r_{ij} < 0 \end{cases}$$

由式 (11) 得,
$$H(z) = \nabla^2 E(z) = \frac{\partial}{\partial z^T} \nabla E(z) = \begin{matrix} T + \frac{1}{2} T \{ I - \text{diag}(\text{sgn}(z_1), \text{sgn}(z_2), \dots, \text{sgn}(z_{n+m})) \} + \frac{1}{2} \{ I - \text{diag}(\text{sgn}(z_1), \text{sgn}(z_2), \dots, \text{sgn}(z_{n+m})) \} \end{matrix}$$

显然,能量函数 $E(z)$ 的 Hess 阵是一个半正定实对称矩阵,并且

$$H(z) = T + T + I \quad (12)$$

式 (12) 的数学意义是,如果点列 $\{z^0, z^1, \dots, z^j, \dots\}$ 中的某点 z^j 落在原规划和对偶规划的可行域中,则 $H(z) = T$, 否则 $H(z) > T$. 同时也表明,当算法未收敛时,可用 $T + T + I$ 的某种测度作为 $H(z)$ 的估计量去研究算法步长. 联想到最优化方法中的牛顿法,可以猜测到,算法步长应是该估计量的倒数形式. 下面的结果说明了这一点.

为了书写方便,采用记号 z^k 表示第 k 时刻 z 的取值. 除非声明,后面符号意同. 令 $g(k) = \nabla E(z^k)$, $H(k, z) = \nabla^2 E(z^k + (z^{k+1} - z^k))$, 且 $0 < \mu < 1$.

定理 9 神经网络 (11) 收敛的充分必要条件是 $I - \frac{1}{2} \mu_k H(k, z) > 0, k = 0, 1, 2, \dots$

证明见附录

推论 神经网络 (11) 收敛的充分必要条件是算法步长 μ_k 满足

$$0 < \mu_k < \frac{2}{\max_k}, k = 0, 1, \dots \quad (13)$$

证明见附录

定理 10 神经网络方程 (11) 收敛的一个充分条件是算法步长

$$\mu = \frac{1}{T + 2} \quad (14)$$

证明见附录

5 算法的数值仿真

定理 6 说明由微分方程 (10) 规定的神经网络的初值解存在唯一性, 定理 7 说明网络收敛到最优解, 定理 10 给出了保证算法收敛的步长公式. 下面给出计算机仿真结果.

为验证算法收敛的一致性, 初始点 $z(0) = [x^T(0), y^T(0)]^T$ 的每个分量在 $(-100, 100)$ 范围内随机选取 (均匀分布), $\mu = 1/L$, 例 1 和例 2 的收敛控制误差分别为 $\nabla E(z^k)$

$2 < 10^{-4}$ 和 10^{-7} .

例 1^[2]

$$\min \begin{cases} 12x_1 + 8x_2 + 16x_3 + 12x_4 \\ 2x_1 + x_2 + 4x_3 \geq 2 \\ 2x_1 + 2x_2 + 4x_4 \geq 3 \\ x_i \geq 0, i = 1, 2, 3, 4 \end{cases}$$

例 2^[4]

$$\min \begin{cases} x_0 = 8x_{11} + 14x_{12} + 12x_{13} + 17x_{14} + 11x_{21} + 9x_{22} + 15x_{23} + 13x_{24} + 12x_{31} + 19x_{32} + 10x_{33} + 6x_{34} + 12x_{41} + 5x_{42} + 13x_{43} + 18x_{44} \\ x_{11} + x_{12} + x_{13} + x_{14} = 20 \\ x_{21} + x_{22} + x_{23} + x_{24} = 10 \\ x_{31} + x_{32} + x_{33} + x_{34} = 10 \\ x_{41} + x_{42} + x_{43} + x_{44} = 15 \\ x_{11} + x_{21} + x_{31} + x_{41} = 15 \\ x_{12} + x_{22} + x_{32} + x_{42} = 20 \\ x_{13} + x_{23} + x_{33} + x_{43} = 10 \\ x_{14} + x_{24} + x_{34} + x_{44} = 10 \\ x_{ij} = 0, i, j = 1, 2, 3, 4 \end{cases}$$

表 1 用学习率公式计算例 1 得到的最优解和理论解. (误差精度控制万分之一)

例 1 的解	理论解	计算解
$\max f(x)$	14	13.99999999
x_1	0.5	0.498526466602
x_2	1	1.001473534990
x_3	0	0.000368382728
x_4	0	0.00000000907

表 2 用学习率公式计算例 2 得到的最优解 (精确到十亿分之一). 计算最优解 435. 理论最优解 435^[4]. (误差精度控制百分之一)

例 2 的解 $X(i, j)$	$i=1$	2	3	4
$i=1$	13.863	0	6.1374	0
2	1.1374	5	3.8626	0
3	0	0	0	10
4	0	15	0	0

计算机仿真结果如表 1 和表 2 所示. 由此可以看出, 本文的理论分析是正确的, 给出的算法步长公式是有效的.

6 结论

本文首先构造了对应于线性规划问题的能量函数, 导出了神经网络微分方程, 简要分析了其性质. 然后, 用奇异函数概念导出了能量函数的 Hess 阵, 分析了能量函数及其梯度的变化规律. 用特征值方法界定了算法步长的取值范围, 给出了用李普希兹常数表示的算法步长公式. 算法数值仿真结果表明, 理论分析是正确的, 算法步长公式是有效的.

附录

引理 1 设 $x = (x_1, x_2, \dots, x_n)^T \in R^n$, 则 $\phi(x) = \frac{1}{2} x^T (x - |x|)$ 是可微不减凸函数.

证 (1) 凸性. 在 $(x) = \frac{1}{2} \{ (x_1^2 - x_1 |x_1|) + (x_2^2 - x_2 |x_2|) + \dots + (x_n^2 - x_n |x_n|) \}$ 中, 每一项 $x_i^2 + x_i |x_i| = 0, x_i \geq 0$
 $f - 2x_i^2, x_i < 0$ 是凸函数, 因此 (x) 是凸函数.

(2) 可微不增性. 因 $x^T |x| = x_1 |x_1| + x_2 |x_2| + \dots + x_n |x_n|$, 所以 $\frac{\partial}{\partial x_i} (x^T |x|) = \frac{\partial}{\partial x_i} (x_i |x_i|) = 2|x_i|$, 故 $(x) = (x - |x|)^T \cdot 0$. 因此, (x) 是可微不增凸函数. 证毕.

引理 2 设 $A_k = (a_{ij}^k)$ 是矩阵序列, 记 $|A_k| = (|a_{ij}^k|), A_{k-1} = A_k - |A_k|$, 则 $|A_{k-1} - A_{k-2}| \leq 2|A_k - A_{k-1}|$.

证 $|A_{k-1} - A_{k-2}| = |(A_k - |A_k|) - (A_{k-1} - |A_{k-1}|)|$
 $|A_{k-1} + (|A_k| - |A_{k-1}|)|$

而 $|A_k| - |A_{k-1}| = (|a_{ij}^k|) - (|a_{ij}^{k-1}|) = (|a_{ij}^k| - |a_{ij}^{k-1}|)$

$$|a_{ij}^k - a_{ij}^{k-1}| = |A_k - A_{k-1}|$$

故 $|A_{k-1} - A_{k-2}| \leq 2|A_k - A_{k-1}|$ 证毕

定理 5 的证明

证 (1) $E_1(z) = \frac{1}{2} z^T T_z^{-2} = \frac{1}{2} z^T T_z$ 是半正定二次型函数, 显然可微, 且 $\nabla^2 E_1(z) = T_z^{-1} > 0$, 因而又是凸函数.

(2) $E_3(z) = \frac{1}{2} z^T z^{-2} = -\frac{1}{4} z^T (z - |z|)$ 由引理 1 知, $E_3(z)$ 是可微凸函数.

(3) $E_2(z) = \frac{1}{2} (z + |z|)^{-2}$, 因 $z + |z|$ 是 z 的线性变换, 由线性变换性质和证 (2) 的结果知 $E_2(z)$ 是可微凸函数.

综合 (1) (2) (3), 即得所证. 证毕

定理 9 的证明

证 充分性. 将能量函数 $E(z^{k+1})$ 在 z^k 处展开为泰勒级数

$$E(z^{k+1}) = E(z^k) + (z^{k+1} - z^k)^T g(k) + \frac{1}{2} (z^{k+1} - z^k)^T \cdot H(k, z) (z^{k+1} - z^k) \quad (a1)$$

将式 (11) 写为 $z^{k+1} = z^k - \mu_k \nabla E(z^k)$ 并代入 (a1) 中, 得

$$E(z^{k+1}) = E(z^k) - \mu_k g^T(k) g(k) + \frac{1}{2} \mu_k^2 g^T(k) H(k, z) g(k)$$

$$= E(z^k) - \mu_k g^T(k) [I - \frac{1}{2} \mu_k H(k, z)] g(k) \quad (a2)$$

当算法不做收敛时, $z^k \notin \{z \in R^{n+m} | g(k) = 0\}$ 即 $g(k)$

0. 此时如果 $I - \frac{1}{2} \mu_k H(k, z) > 0$, 则 $E(z^{k+1}) < E(z^k)$. 又由式 (4) 知 $E(z^k) > 0$. 所以, $E(z^k)$ 是严格单调下降有界序列, 存在极限,

$$\lim_k E(z^{k+1}) = E(z) = 0, z \in \{z \in R^{n+m} | E(z) = 0\}$$

顺便指出, $H(k, z)$ 通常是正定的, 因而 μ_k 是一个正数序列.

必要性. 如果神经网络 (11) 收敛, 则 $E(z^k)$ 必须是单调下降有界序列. 由式 (18) 用反证法立证.

定理 9 推论的证明

证 设 $\mu_1^k, \mu_2^k, \dots, \mu_{n+m}^k$ 是 $H(k)$ 的 $n+m$ 个特征值, $H(k)$ 是半正定的, 令 $H(k) = P \mu_k P^T, \mu_k = \text{diag}(\mu_1^k, \mu_2^k, \dots, \mu_{n+m}^k)$ 则 $E(z^{k+1}) = E(z^k) - \mu_k g^T(k) p(I - \frac{1}{2} \mu_k) p^T g(k) = E(z^k) - \mu_k V^T(k) (I - \frac{1}{2} \mu_k) V(k)$. 这里 $V(k) = p^T g(k)$ 相当于坐标轴旋转变换, 它不改变问题的几何结构. 因此,

$$I - \frac{1}{2} \mu_k H(k) > 0 \Leftrightarrow \text{diag}(1 - \frac{1}{2} \mu_1^k, 1 - \frac{1}{2} \mu_2^k, \dots, 1 -$$

$$\frac{1}{2} \mu_{n+m}^k) > 0$$

令 $\mu_{\max}^k = \max_{0 \leq i \leq n+m} \{\mu_i^k\}$, 则当 $I - \frac{1}{2} \mu_{\max}^k > 0$ 时有 $I - \frac{1}{2} \mu_k H(k) > 0$.

故 $0 < \mu_k < \frac{2}{\mu_{\max}^k}$

定理 10 的证明

证 由定理 6 得, $g(k+1) - g(k) \leq L \cdot z^{k+1} - z^k$. 算法未收敛时, $z^{k+1} > z^k$, 因而

$$\frac{g(k+1) - g(k)}{z^{k+1} - z^k} \leq L$$

又, 由 Hess 阵定义知

$$H(k) = \frac{g(k+1) - g(k)}{z^{k+1} - z^k} \leq \frac{g(k+1) - g(k)}{z^{k+1} - z^k} \leq L$$

根据矩阵理论, 对一切 k 存在 $\mu_{\max}^k \leq H(k)$. 因此, 取

$\mu = \frac{1}{L}$, 必有

$$0 < \mu = \frac{1}{L} \leq \frac{1}{\mu_{\max}^k} < \frac{2}{\mu_{\max}^k}, k = 0, 1, 2, \dots \quad \text{证毕}$$

参考文献:

- [1] Cichocki A & Unbehauen R. Neural networks for solving system of linear equations and related problems [J]. IEEE Trans. on Circuits and systems-I, 1992, 39(2): 124 - 137.
- [2] 陈开周. 最优化计算方法 [M]. 西安: 西安电子科技大学出版社, 151 - 168.
- [3] Xia Y, et al. Recurrent neural networks for solving linear inequalities and equations [J]. IEEE Trans. on Circuits and systems-I, 1999, 46(4): 452 - 462.
- [4] Foulds L R. Optimization techniques, an introduction [M]. New York: Springer-verlay, 1981: 418 - 420.

作者简介:

裴炳南 男, 1956 年生. 郑州大学电子工程系教授, 现在西安电子科技大学师从保铮院士攻读博士学位, 从事雷达信号处理和雷达目标识别方面的研究. e-mail: bnpei@rsp.xidian.edu.cn

保铮 男, 1927 年生. 西安电子科技大学教授, 中国科学院院士, 长期从事雷达信息系统、雷达信号检测与处理和现代信号处理方面教学和研究工作.

