

神经网络系统学习过程初探¹⁾

西广成

(中国科学院自动化研究所,北京)

摘 要

本文给出了神经网络系统学习过程的一个随机模型——马尔柯夫模型,从最大熵原理的观点讨论神经网络系统的学习过程,提出将神经网络系统的学习过程分为两个阶段的想法并给出学习过程的算法。

关键词: 最大熵,神经网络,学习过程。

一、引 言

一类人工神经网络系统可看成是由大量简处理单元组成的大规模并行网络系统,它工作在复杂的环境中,经常不断地对新环境、对已存在环境的新特征进行学习、再学习。神经网络系统接受足够强的输入信息时将发生状态转移,在转移途中,实行状态间的竞争与协调,一旦哪种状态压倒对方而成为优胜者,它就成为学习过程的稳定状态。对应于不同的判决规则,神经网络系统具有不同的稳定状态,这是生物神经网络系统所显现的一个基本特征,试图将一类人工神经网络系统纳入一种数学框架是本文的目的。

二、神经网络系统学习过程的随机模型

设环境 E 由有限个状态 $E_e, e = \{1, 2, \dots, a\}$ 组成,即 $E = \{E_1, E_2, \dots, E_a\}$; 神经网络系统 N 由有限个状态 $N_f, f = \{1, 2, \dots, b\}$ 组成,即 $N = \{N_1, N_2, \dots, N_b\}$. 神经网络系统 N 的输出 $r \in R$ (R 是所有可能输出的集合) 是影响环境 E 状态转移的环境 E 的输入,环境 E 的状态转移是随机的(将时而出现的确定性情况视为概率为 1 的发生事件),仅依赖于目前的状态 E_i 和神经网络系统的输出 $r \in R$. 转移概率 $p_{ij}(r)$ 表示在神经网络系统输出为 r 的情况下环境 E 由状态 E_i 到状态 E_j 的转移概率。

假定制约着神经网络系统的心理的生理的以及其它的各种规则本身及这些规则强度大小均认为不变化,于是不妨认为系统 N 的状态转移仅决定于目前的状态 N_i 和来自环境的影响下一个状态 N_j 和输出 $r \in R$ 的信息. 一般情况下,系统 N 将随机地决定 N_j 及 $r \in R$. 环境 E 和系统 N 的目前状态 (E_i, N_j) 对应着唯一的输出 $r \in R$. 在联合状态 $(E_i,$

本文于 1989 年 6 月 3 日收到。

1) 国家自然科学基金资助课题。

N_i) 的条件下, N 的输出为 r 的概率是 $p(r|(E_s, N_i))$; 系统 N 到达下一个状态 N_j 的概率用 $p_{ij}(E_s)$ 表示.

根据以上分析, 可得联合状态 $\mathbf{c}_u(E_s, N_i)$ 到联合状态 $\mathbf{c}_{u'}(E_t, N_j)$ 的转移概率

$$\begin{aligned} p(\mathbf{c}_{u'} = (E_t, N_j) | \mathbf{c}_u = (E_s, N_i)) \\ = p_{ij}(E_s) \sum_{r \in R} p_{sr}(r) p(r|(E_s, N_i)). \end{aligned} \quad (1)$$

将式(1)称为神经网络系统学习过程的马尔柯夫模型. 此模型详细地描述了环境 E 和系统 N 的相互作用.

下面讨论神经网络系统 N 由状态 $i, i \in f$ 到状态 $j, j \in f, i \neq j$ 的转移就是在这一模型所包含的全部背景下进行的. 但是, 着眼点不在于神经网络系统 N 和环境 E 具体怎样发生作用以及在怎样的相互作用下发生怎样的转移, 着眼点定在对神经网络系统的状态观测实现 N_j 以及这种状态观测实现的转移概率. 这样做, 既简单又不影响实现本文目的.

三、神经网络系统学习过程和最大熵原理

由 q 个二值单元构成的神经网络系统有 $b = Z^q$ 个状态, 每个状态用一个点 $i \in f$ 表示. 从状态 i 到状态 j 的转移用一些带箭头的线相互连接而成, 这样形成的图称为状态迁移图^[4]. 状态迁移图(假定由状态 i 经 m 步到状态 j) 在平面上具有图 1 所示的形式:

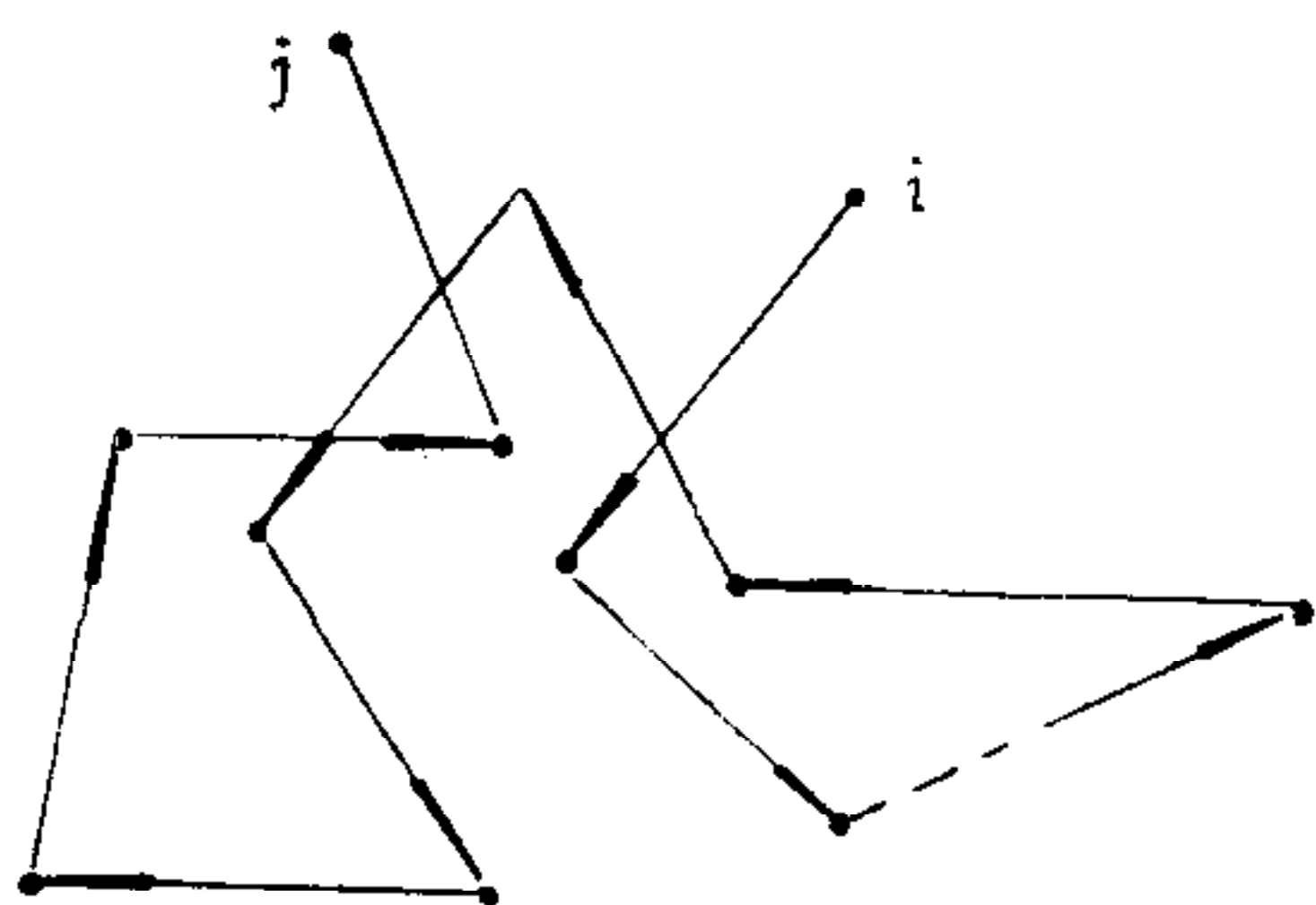


图 1 状态迁移图的一个实例

由状态 i 到状态 j 转移概率的性质由定理 1 给出.

定理 1. 在上述状态迁移图中, 状态 i 经 m 步到达状态 j , 则 $\lim_{m \rightarrow \infty} p_{ij}^{(m)}$ 遵从零均值方差为 $\sigma^2 t$ 的高斯分布. 证明见附录 A.

处在远离平衡状态的系统, 它朝向平衡状态移动比朝向相反方向移动的可能性大得多, 并迟早要到达平衡状态. 设系统 N 由初始状态 i_0 经 $(n-1)$ 个状态到达某平衡状态 i_n , 由定理 1, 可得

$$\begin{aligned} p(i_0, i_1, \dots, i_n) &= p_{i_0 i_1} p_{i_1 i_2} \dots p_{i_{n-1} i_n} \\ &= p_{i_0} \frac{1}{\prod_{i=1}^n \sqrt{2\pi\sigma_i^2 t_i}} \exp\left(-\sum_{i=1}^n \frac{x^2}{2\sigma_i^2 t_i}\right), \end{aligned} \quad (2)$$

求出 p_{i_0} 并代入上式, 有

$$p(i_0, i_1, \dots, i_n) = \frac{\sqrt{\sum_{i=1}^n 2\sigma_i^2 t_i}}{\sqrt{\pi}} \exp\left(-\sum_{i=1}^n \frac{1}{2\sigma_i^2} \frac{x^2}{t_i}\right)$$

$$= \frac{1}{\Phi} \exp\left(-\sum_{i=1}^n \lambda_i f_i(x)\right), \quad (3)$$

其中

$$\begin{cases} \Phi = \frac{\sqrt{\pi}}{\sqrt{\sum_{i=1}^n 2\sigma_i^2 t_i}}, \\ \lambda_i = \frac{1}{2\sigma_i^2}, \\ f_i(x) = \frac{x^2}{t_i}. \end{cases} \quad (4)$$

在处理随机现象时,若已知系统的概率分布,即系统中各个事件发生的概率,则可以计算系统分布的均值、方差,进而考察系统的各种统计性质.但在实际上,很少能准确地知道系统的概率分布.多数情况下,仅知道为确定系统概率分布所需要的部分信息.从符合这部分信息的若干概率分布中,用某种办法选出一种特定分布,根据这个分布进行系统随机现象的统计分析.问题是在给定部分信息的情况下,怎样选出一个最佳分布,即应当根据什么样的规则、原理选择这样的概率分布,这就是最大熵原理所要回答的问题.关于最大熵原理有以下的引理.

引理. 设域 $S \in \mathbf{R}^n$, X 是遵从 S 上连续分布的 n 维随机变数, $f_1(x), f_2(x), \dots, f_n(x)$ 是定义在 S 上的实函数,随机变数 $f_1(X), \dots, f_n(X)$ 的期望是

$$E[f_j(X)] = m_j, \quad j = 1, \dots, n, \quad (5)$$

其中 m_j 满足 $\Phi, \lambda_j, j = 1, 2, \dots, n$ 的方程组:

$$\begin{cases} \Phi = \int_S \exp\left\{-\sum_{i=1}^n \lambda_i f_i(x)\right\} dx, \\ \frac{1}{\Phi} \int_S f_j(x) \exp\left\{-\sum_{i=1}^n \lambda_i f_i(x)\right\} dx = m_j, \end{cases} \quad (6)$$

$$\frac{1}{\Phi} \int_S f_j(x) \exp\left\{-\sum_{i=1}^n \lambda_i f_i(x)\right\} dx = m_j, \quad (7)$$

则此时使得连续熵最大的 X 的分布密度函数为

$$p^*(x) = \frac{1}{\Phi} \exp\left\{-\sum_{i=1}^n \lambda_i f_i(x)\right\}, \quad x \in S, \quad (8)$$

并且是唯一的. 证明见附录 B.

由式(3)和引理可知,神经网络系统在获得足够强的输入信息时开始它的学习过程.该过程从远离平衡状态到达平衡状态,平衡状态即为最大熵状态.将以上讨论结果表述为如下一般性的定理.

定理 2. 设神经网络系统 N 和它所在的环境系统 E 组成的复合系统与外界孤立,神经网络系统的学习过程为由式(1)给出的马尔柯夫过程.则在神经网络系统获得足够强的输入信息时,开始自己的学习过程,该过程从远离平衡状态最终到达平衡状态.平衡状态是系统的最大熵状态,最大熵状态的充分必要条件是系统的分布密度函数,由式(8)给出.

上面讨论的神经网络系统的学习过程是神经网络系统对环境系统进行抽象确立被求

解问题的过程,这是学习过程的第一阶段。下面讨论学习过程的第二阶段,即求解问题的阶段。

假定神经网络系统概率 1 地具有反映环境特征的先验知识。神经网络系统学习环境规律性过程的第二阶段是环境规律性在神经网络系统中引起的某种不确定性的极小化过程。将这种不确定性定义为 Shannon 熵函数,进而推导出学习算法,即神经网络的学习规则。

四、学习算法

由统计力学知识,当神经网络系统获得能量 E_k (假定神经网络系统与环境系统之间交换的物理量是能量)并和环境系统处于平衡状态,式(8)有如下的离散形式:

$$p_{kr}^* = \frac{1}{\Phi} \exp(-\lambda E_k), \quad (k = 1, 2, \dots, l; r = 1, 2, \dots, 6) \quad (9)$$

其中

$$\Phi = \sum_{k=1}^l \sum_{r=1}^6 \exp(-\lambda E_k) = b \sum_{k=1}^l \exp(-\lambda E_k). \quad (10)$$

引入关系式^[2]

$$E_k = - \sum_{b>a} w_{ab} s_a^{kr} s_b^{kr}, \quad (11)$$

将(11)式代入(9)式,得到

$$p_{kr}^* = \frac{1}{\Phi} \exp\left(\lambda \sum_{b>a} w_{ab} s_a^{kr} s_b^{kr}\right). \quad (12)$$

神经网络的不确定性定义为

$$I = - \sum_r p_{kr}^* \ln p_{kr}^*, \quad (13)$$

神经网络系统的学习是通过改变各神经元之间的连接权重 w_{ab} (假定 $w_{ab} = w_{ba}$) 实现的。为实现 I 的梯度下降,需知 I 关于每个权重的偏导数。

定理 3. 假定神经网络系统的单元 a 和 b 之间的连接权重为 w_{ab} , $w_{ab} = w_{ba}$; 神经网络的不确定性 I 如式(13)定义。则

$$\frac{\partial I}{\partial w_{ab}} = \lambda(P_{ab}(k) - P'_{ab}(k)), \quad (14)$$

其中

$$P_{ab}(k) = b \sum_r p_r^* p_{kr}^* s_a^{kr} s_b^{kr}, \quad (15)$$

$$P'_{ab}(k) = \sum_r p_{kr}^* s_a^{kr} s_b^{kr}, \quad (16)$$

$$p_r^* = \sum_{k=1}^l p_{kr}^*, \quad (17)$$

$\lambda = 1/kT$, k 是玻尔兹曼常数, T 是温度, 控制参数. $P_{ab}(k)$ 反映神经网络系统实际学到的“本领”, $P'_{ab}(k)$ 反映神经网络系统应该学到的“本领”(即环境规律性要求神经网络学到的“本领”). s_a^{kr} 等于 1, 只须单元 a 导通, 否则 s_a^{kr} 等于 0. 由定理 3 可得

$$\Delta w_{ab} = \beta(P_{ab}(k) - P'_{ab}(k)). \quad (18)$$

在实际训练神经网络系统过程中, 怎样更好地决定 β 值是有待于进一步探讨的问题.

定理 3 的证明见附录 C.

附 录

A. 定理 1 的证明

将状态迁移图(图 1) 投影到一维空间即直线上, 只考虑不相交的时间区间. 状态迁移的位置每隔 Δt 时间沿纵坐标轴上下跳跃 $+\delta$ 或 $-\delta$, 在此时间间隔上其值为常数. 这样形成的转移轨迹用 $\xi(t)$ 表示. 显然有

$$\Delta \xi(m) = \xi(m\Delta t) - \xi((m-1)\Delta t), \quad (A.1)$$

$\Delta \xi(m)$ 是独立随机变量序列. 设 $\xi(0) = 0$, $\xi(m\Delta t) = \xi(m)$, 则有

$$\xi(m) = \sum_{k=1}^m \Delta \xi(k), \quad (A.2)$$

由中心极限定理^[3], 有

$$\lim_{m \rightarrow \infty} P\left(\frac{\xi(m)}{\delta \sqrt{m}} \leq n\right) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^n \exp\left(-\frac{y^2}{2}\right) dy. \quad (A.3)$$

现考虑 $\Delta t \rightarrow 0$ 时的极限性质. $t = m\Delta t$ 时 $\xi(t)$ 的均值和方差分别为

$$E\{\xi(t)\} = 0, \quad (A.4)$$

$$D\{\xi(t)\} = t \frac{\delta^2}{\Delta t}, \quad (A.5)$$

考虑固定时刻 t , 让 $\Delta t \rightarrow 0$, 有

$$\lim_{\Delta t \rightarrow 0} m = \lim_{\Delta t \rightarrow 0} \left[\frac{t}{\Delta t} \right] = \infty.$$

(记号 $[\cdot]$ 表示对“ \cdot ”取整)从式 (A.5) 可看出, 为使在时刻 t , $\Delta t \rightarrow 0$ 时的 $\xi(t)$ 的方差存在且非 0, 必需使 δ^2 与 Δt 为同阶无穷小. 设

$$\delta^2 = \sigma^2 \Delta t, \quad (A.6)$$

令

$$F(t) \triangleq \lim_{\Delta t \rightarrow 0} \xi(t), \quad (A.7)$$

$t = m\Delta t$ 时 $F(t)$ 小于等于 x 的概率为

$$\begin{aligned} P(F(t) \leq x) &= \lim_{\Delta t \rightarrow 0} P(\xi(t) \leq x) = \lim_{\Delta t \rightarrow 0} P\left(\frac{\xi(m\Delta t)}{\delta \sqrt{m}} \leq \frac{x}{\delta \sqrt{m}}\right) \\ &= \lim_{m \rightarrow \infty} P\left(\frac{\xi(m)}{\delta \sqrt{m}} \leq \frac{x}{\sigma \sqrt{t}}\right), \end{aligned}$$

应用式 (A.3) 并作变量替换, 可得

$$\begin{aligned} \lim_{m \rightarrow \infty} P \left(\frac{\xi(m)}{\delta \sqrt{m}} \leq \frac{x}{\sigma \sqrt{t}} \right) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\frac{x}{\sigma \sqrt{t}}} \exp\left(-\frac{y^2}{2}\right) dy \\ &= \frac{1}{\sqrt{2\pi t \sigma}} \int_{-\infty}^x \exp\left(-\frac{u^2}{2\sigma^2 t}\right) du. \end{aligned} \quad (A.8)$$

因此,可得到

$$P(F(t) \leq x) = \frac{1}{\sqrt{2\pi t \sigma}} \int_{-\infty}^x \exp\left(-\frac{u^2}{2\sigma^2 t}\right) du. \quad (A.9)$$

定理 1 证完.

B. 引理的证明

当 X 的分布密度函数为 $p(x)$ 时,式 (5) 可写成

$$\int_S f_i(x) p(x) dx = m_i, \quad x \in S, \quad (A.10)$$

设满足式 (A.10) 条件的 $p(x)$ 的全体为 \mathcal{F} , 从式 (7) 显然可得 $p^*(x) \in \mathcal{F}$. 对任意 $p(x) \in \mathcal{F}$, 由连续熵的性质

$$\begin{aligned} h(x) = h(p(x)) &\leq - \int p(x) \ln p^*(x) dx = - \int_{R^n} p(x) \\ &\times \left(-\ln \Phi - \sum_{i=1}^n \lambda_i f_i(x) \right) dx = \ln \Phi + \sum_{i=1}^n \lambda_i m_i = h(p^*(x)). \end{aligned}$$

唯一性显然. 引理证完. 在离散情况下,引理亦正确.

C. 定理 3 的证明

$$I = - \sum_r p_{kr}^* \ln p_{kr}^* = \sum_r \frac{1}{\Phi} e^{\lambda \sum_{b>a} w_{ab} s_a^{kr} s_b^{kr}} \ln \Phi - \lambda \sum_{b>a} w_{ab} s_a^{kr} s_b^{kr} \sum_r \frac{1}{\Phi} e^{\lambda \sum_{b>a} w_{ab} s_a^{kr} s_b^{kr}}$$

在求导运算中,注意使用公式

$$\Phi = \sum_{k=1}^l \sum_{r=1}^b \exp(-\lambda E_k) = b \sum_{k=1}^l e^{\lambda \sum_{b>a} w_{ab} s_a^{kr} s_b^{kr}}$$

通过简单运算,可得式 (14).

参 考 文 献

- [1] 甘利俊一, 神经网络の数理, 产业図書, 1978.
- [2] Ackley, D. H., Hinton, G. E., Sejnowski, T. J., A Learning Algorithm for Boltzmann Machines, *Cog. Sci.*, **9**(1985), 147—169.
- [3] 王梓坤, 概率论基础及其应用, 科学出版社, 1979, 北京.

A TENTATIVE INVESTIGATION OF THE LEARNING PROCESS OF A NEURAL NETWORK SYSTEM

Xi Guangcheng

(Institute of Automation, Chinese Academy of Sciences)

ABSTRACT

In this paper, a Markov process model for the learning process in a neural network is given. A two stage-procedure and related algorithm for the learning process is suggested.

Key words : Maximum entropy; neural network; learning process.