

Mathematical Neuroscience — An Overview

Shun-ichi Amari

RIKEN Brain Science Institute

Abstract

Mathematical neuroscience studies the fundamental principles of information processing that has realized in the brain through long history of evolution. We use simple and tractable models of neural networks and establish rigorous mathematical theories to study the performances of the models. This reveals the secrets why the brain works well by using a large number of stochastically fluctuating neurons and parallel dynamics of mutual interactions. We give a number of examples of mathematical neuroscience approaches, including statistical neurodynamics, associative memory, neural field theory and dynamics of learning in a hierarchical system including singularity.

1 Introduction

The brain is a very complex system consisting of a huge number of neurons. It forms a system composed of complicated networks, storing information as memories and processing information to make decision. One may say that it is a most sophisticated system that the nature has ever created. There are many approaches of research to the brain.

Brain science searches for the mechanisms and functions of the brain. Researchers study, from the microscopic point of view, the role of genes and proteins to understand the structure and function of neurons by using the techniques of molecular biology. They also use molecular imaging techniques. They search for the molecular mechanisms of synapses and their plastic changes. Systems neuroscientists observe the activities of neurons of behaving animals to study how information is represented in the brain and how it is processed.

Recent developments of measurement technology makes it possible to observe the activities of a human brain by using fMRI, PET, NIRS, EEG, MEG and even multi-electrodes. Cognitive neuroscience makes use of these new technologies. Computational neuroscience builds models of neural networks and studies their behaviors by theory or computer simulations. Researchers try to build realistic models by taking experimental details into account, so that they can explain the behaviors of actual networks by theory. See, for example, Dayan and Abbott, 2001.

The brain is an organ which a living animal is equipped with for information processing. A living object is material, but is special in the sense that it has an ability of self-reproduction. It has emerged through the long history of evolution. Information, expressed in terms of DNA molecules, is necessary for this. Hence, a living object keeps information in material. Further, the brain processes information in the environment and decides actions. Higher-order information processing ability as well as structured memory are necessary for this purpose.

The brain has been developed through long years of evolution, and it processes information in neuronal networks. Material has become to be able to keep information and to process it. From this point of view, one may say that information utilizes material to develop itself. The brain sits at the cross point of material and information. It is important to study the material basis of the brain. At the same time, it is necessary to study the principles of information which have realized in material.

Computational neuroscience builds realistic models as faithful to experimental details as possible. This is necessary for explaining actual processes taking place in the brain. The results can then be compared with experimental findings. This is the orthodox way of theory.

However, there would be another approach. The principles of brain information processing are obviously very different from those of computers. Information is distributed over a huge number of neurons and processed by their dynamic interactions. We need to understand the principles which guarantee efficient, reliable and swift information processing in such a distributed system with dynamic interactions. The nature has found such principles through long years of evolution and the brain is a realization of the principles.

We need abstract models of neural networks in order to reveal the principles of parallel and

dynamic information processing, which are simple enough to be able to analyze their capabilities deeply. Mathematical neuroscience is an area that studies the possibilities, limitations and performances of information processing by using distributed parallel dynamics (Amari, 1990). It aims at establishing mathematical theories, where models are as simple as possible so long as they do not miss essential features. The nature has realized these principles by using realistic neurons, so that the brain is far more complex than the mathematical models, but they share common principles. If we understand such principles, we can then proceed further to search for realistic models and to find evidences of realization of such principles in the brain.

The present paper shows examples of mathematical neurosciences and its methods briefly.

2 Statistical Neurodynamics (1): Simple Random Networks

Let us consider a very simple model of mathematical neuron, called the McCulloch-Pitts neuron. Let us further consider a network composed of n neurons, and let $x_i(t), \dots, x_n(t)$ be the outputs of neurons at discrete time t . We assume that their values are 1 or -1 , representing firing or non-firing of the neurons. Neuron i receives inputs from other neurons, which are $x_1(t-1), \dots, x_n(t-1)$. Let w_{ij} be the connection weight from neuron j to neuron i . We neglect the threshold term for simplicity and also external inputs. Then, the state of each neuron changes over time as

$$x_i(t) = \text{sgn} \left\{ \sum_j w_{ij} x_j(t-1) \right\}, \quad (1)$$

where sgn is the signature function,

$$\text{sgn}(u) = \begin{cases} 1, & u \geq 0, \\ -1, & u < 0. \end{cases} \quad (2)$$

Let us use abbreviated notations: We use vector

$$\mathbf{x}(t) = [x_1(t) \dots x_n(t)] \quad (3)$$

for the state of the network at time t , matrix $W = (w_{ij})$ for connections, and T_W for state transition operator which is nonlinear, so that (1) is rewritten as

$$\mathbf{x}(t) = T_W \mathbf{x}(t-1). \quad (4)$$

We study the dynamics (4) described by the state transition operator T_W .

The properties of T_W obviously depends on the matrix W . Statistical neurodynamics (Amari, 1971; Amari, 1974; Amari, Yoshida and Kanatani, 1977) assumes that W is randomly generated from a probability distribution $p(W)$, and searches for the properties that hold for almost all networks having W generated subject to $p(W)$ when n is large. We give two examples.

When w_{ij} are iid (identical and independently distributed) random variables subject to a distribution $p(w)$, it is easy to apply the central limit theorem and the law of large numbers, when n is large. Here, we assume that $p(w)$ is the Gaussian distribution with mean \bar{w}/\sqrt{n} and variance σ_w^2 . We can easily treat a general distribution by a similar method. What are common aspects of dynamics of these networks? To answer this question, we introduce macroscopic variables. The first one is the activity level (Amari, 1971),

$$X(\mathbf{x}) = \frac{1}{n} \sum x_i. \quad (5)$$

Then, the macroscopic variable at time t is

$$X(t) = \frac{1}{n} \sum x_i(t) \quad (6)$$

and we have a macroscopic state transition law of X , if we could find a function F which satisfies

$$X(t) = F \{X(t-1)\}. \quad (7)$$

We search for such a function F . To this end, we use the fact that, for a given non-random x_1, \dots, x_n ,

$$u_i = \sum_j w_{ij} x_j \quad (8)$$

are independently and identically distributed Gaussian random variables with mean and variance,

$$E [u_i] = \sqrt{n} \bar{w} X, \quad V [u_i] = n \sigma^2. \quad (9)$$

Hence, by the law of large numbers, we have

$$F(X) = \text{Prob} \{u_i > 0\} - \text{Prob} \{u_i < 0\}. \quad (10)$$

This is calculated as

$$F(X) = \Phi\left(\frac{\bar{w}X}{\sigma}\right), \quad (11)$$

where

$$\Phi(u) = \frac{1}{\sqrt{2\pi}} \int_{-u}^u \exp\left\{-\frac{t^2}{2}\right\} dt. \quad (12)$$

It is easy to prove that, given $\mathbf{x}(0)$ at time 0 of which macroscopic variable $X(0)$ is given by

$$X(0) = \frac{1}{n} \sum x_i(0), \quad (13)$$

the next macroscopic variable $X(1)$ is

$$X(1) = \frac{1}{n} \sum x_i(1) = F\{X(0)\} \quad (14)$$

with a probability close to 1 for almost all randomly generated networks as n tends to infinity. However, there is difficulty for proving that

$$X(2) = F\{X(1)\}, \quad (15)$$

or more generally

$$X(t) = F\{X(t-1)\}. \quad (16)$$

This is because, even we show, for a fixed \mathbf{x} ,

$$X' = X(T_W \mathbf{x}) = F(X(\mathbf{x})), \quad (17)$$

this does not automatically guarantees that

$$X'' = X(T_W T_W \mathbf{x}) = F\{F(X)\}. \quad (18)$$

In order to derive (15) or (17), we used the fact that w_{ij} are independent. However when we calculate (18) by using $T_W T_W \mathbf{x}$, the two operators are no more independent. In other words, when we evaluate

$$v_i(t) = \sum w_{ij} x_j(t), \quad (19)$$

$x_j(t)$ are no more constants but depend on W . Therefore, when we deal with a macroscopic dynamics of a random networks, we need to verify if the macroscopic equation of type (18) is

valid or not. Computer stimulants show that the marcoscopie dynamics (16) holds. This is discussed in Rozonoer (1969, A, B, C) and Amari et al. Amari (1974), Amari et al (1977). Its theoretical justification is still an unsolved problem.

We can extend the macroscopic dynamics to a network consisting of a number of subnetworks, each of which consists of random networks of different types of neurons. The macroscopic variables are then $\mathbf{X} = (X_1, \dots, X_m)$, each X_α , $\alpha = 1, \dots, m$, denoting the macroscopic activity of subnetwork α . We have the following macroscopic equation

$$X_\alpha(t+1) = \Phi \left(\frac{\sum \bar{w}_{\alpha\beta} X_\beta}{\sigma_\alpha} \right). \quad (20)$$

We can prove that these dynamical equations show monostable, multistable, oscillatory and chaotic behaviors depending on the probability distributions of $w_{ij}^{(\alpha)}$ within a subnetwork and of $w_{ij}^{(\alpha\beta)}$ across subnetworks (Amari, 1971, 1972A).

It is interesting to search for other types of macroscopic variables. Let us study the stability of the state transition T_W of a simple random network. Assume that microscopic state \mathbf{x} is mapped to \mathbf{y} by T_W

$$\mathbf{y}' = T_W \mathbf{x}. \quad (21)$$

When \mathbf{x} is disturbed and changes to \mathbf{x}' , how much change is caused in the resultant output $\mathbf{y}' = T_W \mathbf{x}'$? This is the problem of stability of T_W . We introduce the Hamming distance of two microscopic states

$$D(\mathbf{x}, \mathbf{x}') = \frac{1}{2n} \sum |x_i - x'_i| \quad (22)$$

to show their difference. The problem is to find the relation

$$D' = K(D) \quad (23)$$

where

$$D' = D(T_W \mathbf{x}, T_W \mathbf{x}'). \quad (24)$$

We can obtain this relation explicitly when $\bar{w} = 0$ and $\sigma^2 = 1$. The answer given by Amari (1974) is

$$K(D) = \frac{2}{\pi} \sin^{-1} \sqrt{D} \quad (25)$$

The state transition of a random network has a very specific property: When D is small,

$$K(D) \approx \frac{2}{\pi} \sqrt{D}. \quad (26)$$

Since the derivative of $K(D)$ at $D = 0$ is infinity, this means that a small deviation is enlarged extraordinary. So this type of network is convenient for checking small deviations of inputs, since the difference is expanded in the outputs.

If the macroscopic dynamics

$$D_{t+1} = K(D_t) \quad (27)$$

holds for two serieses $\mathbf{x}_t = T_W^t \mathbf{x}_0$ and $\mathbf{x}'_t = T_W^t \mathbf{x}'_0$ and

$$D_t = D(\mathbf{x}_t, \mathbf{x}'_t), \quad (28)$$

we can study the state transition diagram of a random network. Since the number of states is 2^n , we have a state-transition graph having 2^n nodes which correspond to 2^n \mathbf{x} 's. From each node \mathbf{x} , there emerges one branch and its destination is $T_W \mathbf{x}$. Hence, the graph has 2^n directed branches. The state transition graph is highly different from the random state transition graph where the state transition branches are randomly assigned to each node. First, a random neural network has the small-world property, where a small number of nodes monopolize branch outlets. Let p_i be the probability that a node has i ancestors,

$$p_i = \text{Prob} \left\{ \left| T_W^{-1} \mathbf{x} \right| = i \right\}. \quad (29)$$

Then, the distribution is subject to the power law

$$p_i \propto \frac{1}{i^{3-\alpha}}, \quad (30)$$

for some $0 < \alpha < 1$. See Amari (1974). The small-world network has a number of interesting properties, and one of them is that the transient period of dynamics is very short, and the dynamics of

$$\mathbf{x}(t+1) = T_W \mathbf{x}(t) \quad (31)$$

converges quickly to its attractor states.

3 Statistical Neurodynamics (2): Associative Memory

The conventional model of associative memory networks can also be analyzed by statistical neurodynamics. See Amari (1972B), Hopfield (1982) and many others. But it has quite different behaviors from the previous random networks, because the connections w_{ij} are random but not independent. Memorized patterns in this network are represented by the states \mathbf{x} which satisfy

$$T_W \mathbf{x} = \mathbf{x}. \quad (32)$$

Let ξ_1, \dots, ξ_m be m patterns to be memorized in a network such that they are attractors of the dynamics

$$T_W \xi_\alpha = \xi_\alpha, \quad \alpha = 1, \dots, m. \quad (33)$$

Starting at any initial state \mathbf{x}_0 , it is expected that the state reaches one of the memorized pattern ξ_α by the state transition, finding a memorized pattern. When these patterns are generated randomly, we compose the connection matrix in such a way that

$$w_{ij} = \frac{1}{m} \sum \xi_{\alpha i} \xi_{\alpha j} \quad (34)$$

This is a standard way of associative memory. When these patterns are orthogonal, it is easy to see that the patterns are attractors, satisfying

$$T_W \xi_\alpha = \text{sgn} \left(\frac{1}{m} \sum_\beta \xi_\alpha (\xi_\alpha \cdot \xi_\beta) \right) = \xi_\alpha. \quad (35)$$

When ξ'_α 's are randomly generated subject to independent and identical probabilities with $\text{Prob}\{\xi_{\alpha i} = 1\} = \text{Prob}\{\xi_{\alpha i} = -1\} = 1/2$, they are asymptotically orthogonal.

The connection weight matrix W is random, but the component w_{ij} are not independent, because they are composed of (34). The law of stability is interesting in this case. The statistical analysis shows that

$$D' = K(D), \quad K(D) = 1 - \frac{1}{2} \Phi \left\{ -\frac{(1-2D)n}{m} \right\}. \quad (36)$$

Therefore, when D is small, $D = K(D)$ is nearly equal to 0. Hence the state transition of recalling a memorized pattern is very stable. It is interesting to show that the state transition

graph of an associative memory model has a fractal structure (Amari and Maginu, 1988). Here the basin of attractor has very complicated strange boundaries.

4 Pattern Dynamics in Neural Field

A neural field is a continuous version of neural networks, where neurons are located in a 2-dimensional field like cortices. A neural field consists of layers including different types of neurons. We give a general equation describing the dynamics of excitation in the field,

$$\frac{\partial u_i(\mathbf{r}, t)}{\partial t} = -u_i(\mathbf{r}, t) + \sum \int w_{ij}(\mathbf{r}, \mathbf{r}') f\{u_j(\mathbf{r}', t)\} d\mathbf{r}' + I_i(\mathbf{r}, t). \quad (37)$$

Here, $\mathbf{r} = (r_1, r_2)$ is coordinates of the field, denoting the locations of neurons, $u_i(\mathbf{r}, t)$ is the average membrane potential of neurons at location \mathbf{r} in the i -th layer at time t ,

$$z_i(\mathbf{r}, t) = f[u_i(\mathbf{r}, t)] \quad (38)$$

is the output of neurons whose average membrane potential is u_i and f is called the activation function. The functions $w_{ij}(\mathbf{r}, \mathbf{r}')$ are connection weights from the neurons at location \mathbf{r}' in the j -th layer to the neurons at location \mathbf{r} in the i -th layer. $I_i(\mathbf{r}, t)$ are external inputs to the neurons of i -th layer at position \mathbf{r} at time t minus threshold.

This types of dynamics was proposed by Willson and Cown (1973). A simplified version of neural field used by Amari (1977) is a one-dimensional field of one layer, where the activation function is the step function

$$f(u) = \begin{cases} 1, & u \geq 0, \\ 0, & u < 0, \end{cases} \quad (39)$$

and the connection weight function is homogeneous and isotropic, having lateral-inhibition type,

$$w(r, r') = w(|r - r'|), \quad (40)$$

that is, $w(r)$ is positive when r is small and negative otherwise.

Given initial stimuli, we study how an initial excitation pattern $u(r, 0)$ develops in the dynamics. When $u(r) > 0$ holds in a region R of the field, neurons in the region are excited,

$z(r) = 1$, $r \in R$. We call such a region an excited region, that is a region

$$R = \{r | u(r) > 0\}. \quad (41)$$

When an excited region is an interval $R = (a_1, a_2)$, it is called a localized excitation pattern or a bump. We first focus on the dynamics of localized excitation patterns.

Let us consider the case where I is constant over the field. In order to analyze how the excited interval develops over time, we remark its boundaries a_1 and a_2 . They change over time so that we write $a_1(t)$ and $a_2(t)$ for the boundaries of the excited interval at time t . The boundaries satisfy

$$u \{a_1(t), t\} = u \{a_2(t), t\} = 0 \quad (42)$$

at t . After a short time dt , the boundaries change and they satisfy

$$u \{a_i(t + dt), t + dt\} = 0, \quad i = 1, 2. \quad (43)$$

We have, by Taylor expansion,

$$\frac{\partial u(a_i, t)}{\partial t} + \frac{\partial u(a_i, t)}{\partial r} \frac{da_i}{dt} = 0. \quad (44)$$

The time derivative of $u(a_i, t)$ satisfies (37), and the integration part is written as

$$\int w(r - r') f \{u(r', t)\} dr' = \int_{a_1}^{a_2} w(r - r') dr', \quad (45)$$

since the activation function is 1 in the excited interval (a_1, a_2) and 0 otherwise. We define the integration of the connection function by

$$W(r) = \int_0^r w(s) ds. \quad (46)$$

Then, we have

$$\int_{a_1}^{a_2} w(a_2 - r') dr' = W(a_2 - a_1). \quad (47)$$

We put the derivatives of the waveform of pattern $u(r, t)$ at the boundaries,

$$c_1 = \frac{\partial u(r_1, t)}{\partial r}, \quad c_2 = -\frac{\partial u(r_2, t)}{\partial t}. \quad (48)$$

Taking these into account, we have the equations describing the dynamics of the boundary points,

$$\frac{dr_i(t)}{dt} = -\frac{1}{c_i} W(r_2 - r_1) + I, \quad i = 1, 2. \quad (49)$$

Let

$$a(t) = a_2(t) - a_1(t), \quad (50)$$

which is the length of the excited interval. It develops as

$$\frac{da(t)}{dt} = \frac{1}{c} \{W(a) + I\}. \quad (51)$$

The equilibrium a is the solution of

$$W(a) + I = 0. \quad (52)$$

Its stability is analyzed (Amari, 1977), and we have the following results:

When

$$W(a) + I = 0 \quad (53)$$

has a solution, it is an equilibrium of the length of an excited region. It is stable when

$$w(a) > 0 \quad (54)$$

is satisfied.

We show that, under a certain condition, we can find a stable equilibrium solution. The field can retain a local excitation pattern without further inputs in such a case. Since the field is homogeneous, such a pattern can exist at any position. Hence the field has an infinite number of attractors (stable equilibrium solutions), which is called a line attractor (Seung, 1988). Assume that an initial stimulus arrives at around position b . The stimulus may be very noisy, distributed around b . Initial excitations are arisen at around b and smoothed by the dynamic interactions in the field. Even after the stimulus disappears, a local excitation pattern remains stably at the position b when the field admits a local excitation pattern. Hence, this mechanism can be used as a working memory. Such a mechanism is believed to be used in memorizing input patterns and processing them. One example is a neural field detecting

orientations of input bars and keeping the orientations for processing them further. The neural field theory explains the mechanism of the tuning curve in the visual cortex.

More general neural fields give richer behaviors (Ermentrout and Cowan, 1979). See review paper by Coombes and Owen, 2005. When a field consists of an excitatory and inhibitory layers, the field admits a stable travelling region. Its velocity can be controlled by stimuli given outside so that such phenomena will be useful for dynamic information processing.

A two-dimensional neural field has much richer properties. See, for example, Coombes and Owen, 2005. First, it possesses a localized excitation pattern or a bump as well. The stability analysis of the two-dimensional bump was given in a Japanese book in 1978 (Amari), and later proved by some others recently. It also has interesting static patterns consisting of a number of bumps or of a ring shape. Their dynamic behaviors are much more interesting. It has an expanding ring pattern, spiral pattern, breathing pattern and others. It is shown recently that it has also moving bumps. The collision of moving bumps shows very interesting dynamical behaviors.

The neural field model is an important component of computational neuroscience. Moreover, there are lots of applications of dynamic neural field theory to psychological phenomena and robot navigation.

5 Dynamics of Learning in Multilayer Perceptron and Gaussian Mixture

We study the dynamical properties of on-line learning in a hierarchical system such as multilayer perceptrons and Gaussian mixture radial basis systems. It is known that the behavior of a student system trained by using examples given by the teacher system converges to the optimal value, although one cannot avoid a local optimum. However, it is also known that the dynamics of such a learning system suffers from the so-called plateau phenomenon, that the trajectory of learning is attracted to a plateau which is not optimal. It takes long time for the system to get rid of such a plateau. The present section analyzes the dynamical behavior of on-line learning in the neighborhood of a plateau, based on Wei et al. 2008, Wei and Amari,

2008, Cousseau, Ozeki and Amari, 2008.

Plateaus are given rise to by singularities of the parameter space on which learning takes place. This is caused by the symmetric structure of a hierarchical system and we cannot avoid them (Amari, Park and Ozeki, 2006).

Let us consider a function

$$y = k(\mathbf{x}), \quad (55)$$

where \mathbf{x} is an n -dimensional vector $\mathbf{x} = (x_1 \cdots x_n)$ and y is a scalar. We assume that \mathbf{x} is an input to a system and y is its output. In order to realize this function, we use a family of functions

$$F = \{f(\mathbf{x}, \boldsymbol{\theta})\} \quad (56)$$

parameterized by

$$\boldsymbol{\theta} = (\theta_1 \cdots \theta_n). \quad (57)$$

Our problem is to estimate $\boldsymbol{\theta}$ from a series of input-output examples $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ which we call the training examples. Here, \mathbf{x}_t is generated subject to a probability distribution $q(\mathbf{x})$ each time independently, and y_t is the output assumed to be generated by

$$y_t = k(\mathbf{x}_t) + n_t, \quad (58)$$

that is, y_t is the output from the true system contaminated by noise n_t subject to the standard Gaussian distribution. We assume it is subject to $N(0, 1)$ without loss of generality, by rescaling y dividing it by the standard deviation of the noise.

When the true function $k(\mathbf{x})$ from which the training examples are generated belongs to F ,

$$k(\mathbf{x}) = f(\mathbf{x}, \boldsymbol{\theta}^*), \quad (59)$$

the problem is to estimate the true value $\boldsymbol{\theta}^*$ from the training examples. When $k(\mathbf{x})$ does not belong to F , we search for the parameter $\boldsymbol{\theta}^*$ by which $k(\mathbf{x})$ is approximated optimally by $f(\mathbf{x}, \boldsymbol{\theta}^*)$. Hence, this is a regression problem for the data D .

We use an on-line learning approach such that an estimator $\boldsymbol{\theta}_t$ at time t is modified to $\boldsymbol{\theta}_{t+1}$ by using a new example \mathbf{x}_t, y_t at each discrete time $t, t = 1, 2, \dots$. This is a method of sequential

estimation or regression, called on-line learning. The rule of modifying the parameter is written as

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{g}(\mathbf{x}_t, y_t, \boldsymbol{\theta}_t) \quad (60)$$

by using an adequate function \mathbf{g} . Here, η is called a learning constant, which may depend on t . We study the properties of dynamics of the above on-line learning method.

We consider the following two types of the parametric families of functions as F . One is multilayer perceptrons, written as

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j=1}^m v_j \varphi \left(\sum_{i=1}^n w_{ji} x_i \right). \quad (61)$$

Here, the system consists of two layers. The first is the hidden layer and the output of the j -th hidden neuron is a nonlinear function of weighted sum of inputs \mathbf{x} ,

$$z_j = \varphi \left(\sum_i w_{ji} x_i \right) = \varphi(\mathbf{w}_j \cdot \mathbf{x}), \quad (62)$$

where φ is a sigmoidal function and we use

$$\varphi(u) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u \exp \left\{ -\frac{1}{2}s^2 \right\} ds, \quad (63)$$

as an example. The second layer consists of a single neuron, called the output neuron and the output is a linear summation of the outputs of the hidden neurons,

$$z = f(\mathbf{x}, \boldsymbol{\theta}) = \sum v_j z_j. \quad (64)$$

The parameters to be estimated are

$$\boldsymbol{\theta} = (v_1, \dots, v_m; w_{11}, \dots, w_{mn}). \quad (65)$$

The second example uses Gaussian radial basis functions given by

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sum v_j \psi(\mathbf{x} - \mathbf{w}_j) \quad (66)$$

where

$$\psi(\mathbf{x}) = \frac{1}{\sqrt{2\pi}^n} \exp \left\{ -\frac{|\mathbf{x}|^2}{2} \right\} \quad (67)$$

is the Gaussian function. This is similar to the first system of multilayer perceptrons. The parameters to be estimated are

$$\boldsymbol{\theta} = (u_1, \dots, v_m ; \mathbf{w}_1, \dots, \mathbf{w}_m). \quad (68)$$

We use the stochastic gradient learning method, minimizing a cost function each time. Here, the cost function $l(\mathbf{x}, y, \boldsymbol{\theta})$ is given by a half of the square of the difference of the actual output y and the hypothetical output $f(\mathbf{x}, \boldsymbol{\theta})$ which the system with parameter $\boldsymbol{\theta}$ emits for given input \mathbf{x} . This is easily calculated for given parameter $\boldsymbol{\theta}$ when an input-output example (\mathbf{x}, y) is given,

$$l(\mathbf{x}, y, \boldsymbol{\theta}) = \frac{1}{2} \{y - f(\mathbf{x}, \boldsymbol{\theta})\}^2. \quad (69)$$

The difference is called the error function,

$$e(\mathbf{x}, y, \boldsymbol{\theta}) = y - f(\mathbf{x}, \boldsymbol{\theta}). \quad (70)$$

The true system of emitting y can be described by a conditional probability distribution of emitting y from $f(\mathbf{x}, \boldsymbol{\theta})$ with additive noise n (58). Then, the conditional probability distribution of y under the condition that input \mathbf{x} is given is

$$p(y|\mathbf{x}, \boldsymbol{\theta}) = \frac{1}{\sqrt{2\pi}} \exp \left[-\frac{1}{2} \{y - f(\mathbf{x}, \boldsymbol{\theta})\}^2 \right] \quad (71)$$

when the parameter is $\boldsymbol{\theta}$. The cost function (69) is interpreted as the negative of log likelihood

$$l(\mathbf{x}, y, \boldsymbol{\theta}) = -\log p(y|\mathbf{x}, \boldsymbol{\theta}). \quad (72)$$

Hence, minimizing the cost function is equivalent to maximizing the log likelihood. This gives the maximum likelihood estimator.

The dynamics of the stochastic gradient on-line learning method is written as

$$\mathbf{g}(\mathbf{x}, y, \boldsymbol{\theta}) = \nabla l(\mathbf{x}, y, \boldsymbol{\theta}) = -e(\mathbf{x}, y, \boldsymbol{\theta}) \nabla f(\mathbf{x}, y, \boldsymbol{\theta}) \quad (73)$$

where ∇ is the gradient with respect to $\boldsymbol{\theta}$,

$$\nabla = \left(\frac{\partial}{\partial \theta_1} \dots \frac{\partial}{\partial \theta_n} \right). \quad (74)$$

In order to analyze the behavior of dynamics of on-line learning, we use the average learning equation

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \langle \nabla l(\mathbf{x}, y, \boldsymbol{\theta}) \rangle \quad (75)$$

where $\langle \cdot \rangle$ denotes the average over all possible input \mathbf{x} by using the probability distribution (72). We further use the continuous time approximation, and analyze the continuous time average learning equation,

$$\frac{d\boldsymbol{\theta}(t)}{dt} = - \langle \nabla l(\mathbf{x}, y, \boldsymbol{\theta}) \rangle. \quad (76)$$

Let $S = \{\boldsymbol{\theta}\}$ be the parameter space. When

$$f(\mathbf{x}, \boldsymbol{\theta}_1) \neq f(\mathbf{x}, \boldsymbol{\theta}_2) \text{ for } \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2 \quad (77)$$

holds, the system is identifiable, because only one $\boldsymbol{\theta}$ corresponds to a function $f(\mathbf{x}, \boldsymbol{\theta})$. However, when there exist at least two points for which

$$f(\mathbf{x}, \boldsymbol{\theta}_1) = f(\mathbf{x}, \boldsymbol{\theta}_2), \quad \boldsymbol{\theta}_1 \neq \boldsymbol{\theta}_2 \quad (78)$$

holds, the system is not identifiable. Let $R(\boldsymbol{\theta})$ be a subset of S in which all the output function is identical,

$$f(\mathbf{x}, \boldsymbol{\theta}) = f(\mathbf{x}, \boldsymbol{\theta}'), \quad \boldsymbol{\theta}' \in R(\boldsymbol{\theta}). \quad (79)$$

Then, the all parameter points $\boldsymbol{\theta}'$ in $R(\boldsymbol{\theta})$ are equivalent from the point of view of the input-output behavior. There are many such subsets $R(\boldsymbol{\theta})$'s. We divide the space S by this equivalence relation. Then all the points in a $R(\boldsymbol{\theta})$ reduces to one behavior point. The topological structure of S is destroyed by this division, because every $R(\boldsymbol{\theta})$ reduces to a single point. This gives rise to singularities. We call $R(\boldsymbol{\theta})$ a singular region. We give two typical examples of singularity in a hierarchical system.

1) Eliminating singularity:

When $v_i = 0$, whatever \mathbf{w}_i is, $v_i \varphi(\mathbf{w}_i \cdot \mathbf{x}) = 0$, so that $R = \{\boldsymbol{\theta} | v_i = 0, \mathbf{w}_i \text{ arbitrary}\}$ is a singular region. Since neuron i is eliminated because of $v_i = 0$, it is called an eliminating singularity region.

2) Overlapping singularity:

When $\mathbf{w}_i = \mathbf{w}_j = \mathbf{w}$, so long as $v_i + v_j = v$ holds,

$$v_i\varphi(\mathbf{w}_i \cdot \mathbf{x}) + v_j\varphi(\mathbf{w}_j \cdot \mathbf{x}) = v\varphi(\mathbf{w} \cdot \mathbf{x}). \quad (80)$$

So we cannot identify each of v_i and v_j . The region

$$R = \{\boldsymbol{\theta} \mid \mathbf{w}_i = \mathbf{w}_j = \mathbf{w}, v_i + v_j = v\} \quad (81)$$

is a singular region, depending on \mathbf{w} and v . Since the two neurons become identical in this case, we call it an overlapping singularity.

There are other singularities where three neurons overlap and so on, but the above two are typical and other singularities are intersections of the above singularity regions. Therefore, we treat a simple system consisting of two hidden neurons, where $\boldsymbol{\theta} = (v_1, v_2, \mathbf{w}_1, \mathbf{w}_2)$ and study the dynamical behavior of learning in such a system.

In the critical region $R(v, \mathbf{w})$, only one neuron suffices, where $\mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}$ or one of v_1 and v_2 is 0. In this case, its input-output function is written by using a single neuron,

$$f(\mathbf{x}, \boldsymbol{\theta}) = v\varphi(\mathbf{w} \cdot \mathbf{x}). \quad (82)$$

The situation is the same in the Gaussian mixture case.

We introduce a new coordinate system $\boldsymbol{\xi} = (\mathbf{v}, w, \mathbf{u}, z)$, defined by

$$\begin{aligned} \mathbf{u} &= \mathbf{w}_2 - \mathbf{w}_1, & \mathbf{w} &= \frac{1}{2}(\mathbf{w}_1 + \mathbf{w}_2), \\ z &= \frac{v_1 - v_2}{v_1 + v_2}, & v &= v_1 + v_2. \end{aligned} \quad (83)$$

Then, the singular region is represented by $\mathbf{u} = 0$ or $z = \pm 1$. Given (v, \mathbf{w}) , singular region $R(v, \mathbf{w})$ exists. A singular region \mathcal{R} consists of three parts, one given by $\mathbf{u} = 0$, that is, $\mathbf{w}_1 = \mathbf{w}_2$, which we call the overlapping singular region, and the other is $z = \pm 1$, that is $w_1 = 0$ or $w_2 = 0$, in which one neuron is eliminated. We have $R = R_1 \cup R_2$,

$$R_1(\mathbf{v}, w) = \{\boldsymbol{\xi} \mid \mathbf{u} = 0, z : \text{arbitrary}; \mathbf{v} \text{ and } w \text{ fixed}\}, \quad (84)$$

$$R_2(\mathbf{v}, w) = \{\boldsymbol{\xi} \mid \mathbf{u} : \text{arbitrary}; z = \pm 1, \mathbf{v} \text{ and } w \text{ fixed}\}. \quad (85)$$

In the subspace in which \mathbf{v} and w are fixed constants, R_1 is a line in which z changes, and R_2 consists of two n -dimensional surfaces. The singular region R is a composite of one line and two surfaces.

We study the dynamical flow near the singularity $R_1(\mathbf{v}, w)$, by using the Taylor expansion with respect to \mathbf{u} , where $|\mathbf{u}|$ is small. The average learning equation is expanded as

$$\dot{\mathbf{u}} = \eta \frac{w}{2} (1 - z^2) \langle e\varphi''(\mathbf{v} \cdot \mathbf{x})(\mathbf{u} \cdot \mathbf{x})\mathbf{x} \rangle + O(|\mathbf{u}|^2), \quad (86)$$

$$\dot{z} = -\frac{1+z^2}{w} \eta \langle \varphi'(\mathbf{v} \cdot \mathbf{x})(\mathbf{u} \cdot \mathbf{x}) \rangle + O(|\mathbf{u}|^2). \quad (87)$$

We can integrate (86) and (87) in the subspace in which S^* where $\mathbf{v} = \mathbf{v}^*$ and $w = v^*$ are the optimal values, that is $v^* \varphi(\mathbf{w}^* \cdot \mathbf{x})$ is the optimal approximation of the target foundation $k(\mathbf{x})$ by using only a single neuron. This gives the flows or trajectories of dynamics in the neighborhood of $\mathbf{u} = 0$ of S^* . They are given by

$$\frac{1}{2} \mathbf{u}^T \cdot \mathbf{u} = \frac{2w^{*2}}{3} \log \left[\frac{(z^2 + 3)^2}{|z|} \right] + c \quad (88)$$

where c is an arbitrary constant.

6 Milnor attractor

The line $R_1(\mathbf{v}^*, w^*)$ in S^* is a critical line of the dynamics of learning. Under a certain condition, we show that the line is divided into two parts such that one part is stable (attractive) and the other is unstable (repulsive). In such a case, there are trajectories attracted to the stable part, and the basin of attraction has a finite measure. This is the Milnor type attractor (Milnor, 1985), different from a saddle point. More strongly, the basin of attraction includes a neighborhood of this part, so that all the points near this part are once attracted to it. After being attracted to it, the parameters still move randomly in $R_1(\mathbf{v}^*, w^*)$ by stochastic fluctuations until they reach the unstable part (Fukumizu and Amari, 2000). Then, the parameters leave R_1 eventually.

We can prove that R_1^* is unstable (saddle) when $A = (A_{ij})$, $A_{ij} = \langle e\phi_{ij} \rangle$, where

$$\phi_{ij} = \frac{\partial^2}{\partial w_i \partial w_j} \varphi(\mathbf{w} \cdot \mathbf{x}) \quad (89)$$

includes both positive and negative eigenvalues. When $w^* A$ is positive-semi-definite, the interval $1 < z^2$ is attractive, and the remaining part $1 > z^2$ is repulsive. Hence, the interval

$1 < z^2$ behaves like a Milnor attractor. This is the plateau phenomenon, taking long time before getting rid of it. When w^*A is negative-semi-definite, the part $1 > z^2$ is attractive and the other part is repulsive.

7 Conclusions

We have shown examples of mathematical neuroscience approaches. They include 1) the macro-dynamics of randomly connected networks and the characteristics of state transition are elucidated. 2) The dynamics of associative memory networks is elucidated from the point of view of macro-dynamics. 3) Pattern formations and pattern interactions in a neural field are investigated by using integro-differential equations. 4) It is studied how the singularity existing in a hierarchical system affects dynamics of learning.

Mathematical approaches are necessary for finding the principles of information processing which the brain utilizes. We need to use simple and tractable models for carrying mathematical analysis. The models are different from the real biological neural networks, but we hope the mathematical theory clarifies the principles from which we can construct realistic models. Mathematical neuroscience has not been fully established and much more efforts should be devoted for constructing such theories.

References

- [1] P. Dayan and L.F. Abbott, *Theoretical Neuroscience: Computational and Mathematical Modeling of Neural Systems*, MIT Press, 2001.
- [2] S. Amari, Mathematical Foundations of Neurocomputing, *Proceedings of the IEEE*, 78, 9, 1443–1463, 1990.
- [3] S. Amari, Characteristics of Randomly Connected Threshold-Element Networks and Network Systems, *Proc. IEEE.*, 59, 1, 35–47, 1971.
- [4] S. Amari, A Method of Statistical Neurodynamics, *Kybernetik*, 14, 201–215, (Heft 4) 1974.

- [5] S. Amari, K. Yoshida and K. Kanatani, A Mathematical Foundation for Statistical Neurodynamics, *SIAM J. Appl. Math.*, 33, 95–126, 1977.
- [6] L.I. Rozonoer, Random logical nets; I, *Avtomat. Telemekh.*, 5, 137–147, 1969A.
- [7] L.I. Rozonoer, Random logical nets; II, *Avtomat. Telemekh.*, 6, 99–109, 1969B.
- [8] L. I. Rozonoer, Random logical nets; III. *Avtomat. Telemekh.*, 7, 127–136, 1969C.
- [9] S. Amari, Characteristics of Random Nets of Analog Neuron-Like Elements, *IEEE Trans. Systems, Man and Cybernetics*, SMC-2, 5, 643–657, November 1972A.
- [10] S. Amari, Learning Patterns and Pattern Sequences by Self-Organizing Nets of Threshold Elements, *IEEE Trans. Computers*, C-21, 11, 1197–1206, 1972B.
- [11] J.J. Hopfield, Neural network and physical systems with emergent collective computational abilities, *Proceedings of the National Academy of Sciences of the United States of America*, 79, 8, 2554–8, 1982
- [12] S. Amari and K. Maginu, Statistical neurodynamics of associative memory, *Neural Networks*, 1, 1, 63–73, 1988.
- [13] H.R. Wilson and J.D. Cowan, A mathematical theory of the functional dynamics of cortical and thalamic nervous tissue, *Kybernetik*, 13, 55–80, 1973.
- [14] S. Amari, Dynamics of Pattern Formation in Lateral-Inhibition Type Neural Fields, *Biological Cybernetics*, 27, 77–87, 1977.
- [15] H. S. Seung, Continuous attractors and oculomotor control, *Neural Networks*, 11, 1253–58, 1998.
- [16] G.B. Ermentrout and J.D. Cowan, A mathematical theory of visual hallucinations patterns, *Biological Cybernetics*, 34, 137–150, 1979.
- [17] S. Coombes and M.R. Owen, Bumps, breathers, and waves in a neural network with spike frequency adaptation, *Physical Review Letters*, 94, 148102, 2005.

- [18] H. Wei, J. Zhang, F. Cousseau, T. Ozeki and S. Amari, Dynamics of Learning Near Singularities in Layered Networks, *Neural Computation*, 20, 813–843, 2008.
- [19] H. Wei and S. Amari, Dynamics of learning near singularities in radial basis function networks, *Neural Networks*, 21, 989–1005, 2008.
- [20] F. Cousseau, T. Ozeki and S. Amari, Dynamics of Learning in Multilayer Perceptrons Near Singularities, *IEEE Transactions on Neural Networks*, 19, 8, 1313–1328, 2008.
- [21] S. Amari, H. Park and T. Ozeki, Singularities Affect Dynamics of Learning in Neuromanifolds, *Neural Computation*, 18, 1007–1065, 2006.
- [22] J. Milnor, On the concept of attractor, *Commun. Math. Phys.*, 99, 177–195, 1985.
- [23] K. Fukumizu and S. Amari, Local minima and plateaus in hierarchical structures of multilayer perceptrons, *Neural Networks*, 13, 317–327, 2000.