

Information Geometry of α -Projection in Mean Field Approximation

Shun-ichi Amari*, Shiro Ikeda[†] and Hidetoshi Shimokawa[‡]

Abstract

Information geometry is applied to mean field approximation for elucidating its properties in the spin glass model or the Boltzmann machine. The α -divergence is used for approximation, where α -geodesic projection plays an important role. The naive mean field approximation and TAP approximation are studied from the point of view of information geometry, which treats the intrinsic geometric structures of a family of probability distributions. The bifurcation of the α -projection is studied, at which the uniqueness of the α -approximation is broken.

1 Introduction

Mean field approximation uses a simple tractable family of probability distributions to calculate quantities related to a complex probability distribution including mutual interactions. Information geometry, on the other hand, studies intrinsic geometrical structure existing in the manifold of probability distributions (Chentsov, 1982; Amari, 1985; Amari and Nagaoka, 2000). It was used for analyzing performances of learning in the Boltzmann machine (Amari, Kurata and Nagaoka, 1992), the EM algorithm (Amari, 1995), multilayer perceptrons (Amari, 1998), etc. A number of new works appeared recently

*RIKEN, BSI

[†]PRESTO, JST

[‡]Univ. of Tokyo

which treated mean field approximation from the point of view of information geometry (Tanaka [this volume], Tanaka [2000], Kappen [this volume], and Bhattacharyya and Keerthi [1999, 2000]).

The present article studies the relation between mean field approximation and information geometry in more detail. We treat a simple spin model like the SK model or the Boltzmann machine, and study how the mean values of spins can be approximated. It is known (Tanaka, this volume) that, given a probability distribution including complex mutual interactions of spins or neurons, the mean value of each spin is kept constant when it is projected by the m -geodesic to the submanifold consisting of independent probability distributions. On the other hand, the m -projection is computationally intractable for a large system. Instead, its projection by the e -geodesic is easy to calculate. However, the mean value is changed by this, so that only an approximate value is obtained. This gives the naive mean field approximation.

A family of divergence measures named the α -divergence is defined invariantly in the manifold of probability distributions (Amari, 1985; Amari and Nagaoka, 2000). The $\alpha = -1$ divergence is known as the Kullback-Leibler divergence or cross entropy, $\alpha = 1$ as the reverse of the K-L divergence, and $\alpha = 0$ as the Hellinger distance. This concept is closely related to the Rényi entropy (Rényi, 1961; see also Chernoff, 1952; the f -divergence of Csiszár, 1975). These divergence functions give a unique Riemannian metric to the manifold of probability distributions. It moreover gives a family of invariant affine connections named the α -connections where α - and $-\alpha$ -affine connections are dually coupled to each other with respect to the Riemannian metric (Nagaoka and Amari, 1982; Amari, 1985; Amari and Nagaoka, 2000). The α -geodesic is defined in this context. It should be remarked that the Tsallis entropy (Tsallis, 1988) is closely connected to the α -geometry.

We use the α -geodesic projection to elucidate various mean field approximations. The α -projection is the point in the tractable subspace consisting of independent distributions that minimizes the α -divergence from a given true distribution to the subspace. It is given by the α -geodesic that is orthogonal to the subspace at the α -projected point. This gives

a family of the α -approximations, where $\alpha = -1$ is the distribution giving the true mean value and $\alpha = 1$ is the naive mean approximation. Therefore, it is interesting to know how the α -approximation depends on α . We also study the TAP approximation in this framework.

We can prove that the m-projection ($\alpha = -1$) is unique, while α -projection ($\alpha \neq -1$) is not necessarily so. The e-projection ($\alpha = 1$), that is the naive mean field approximation, is not necessarily uniquely solved. Therefore, it is interesting to see how the α -approximation bifurcates depending on α . We calculate the Hessian of the α -approximation which shows the stability or the local minimality of the projection.

2 Geometry of Mean Field Approximation

Let us consider a system of spins or Boltzmann machine, where $\mathbf{x} = (x_1, \dots, x_n)$; $x_i = \pm 1$ denotes the values of n spins. The equilibrium probability distribution is given by

$$q(\mathbf{x}; W, \mathbf{h}) = \exp \{W \cdot X + \mathbf{h} \cdot \mathbf{x} - \psi_q\}, \quad (1)$$

where $W = (w_{ij})$ and $X = (x_i x_j)$ are symmetric matrices and

$$W \cdot X = \frac{1}{2} \sum w_{ij} x_i x_j \quad (w_{ii} = 0), \quad (2)$$

$$\mathbf{h} \cdot \mathbf{x} = \sum h_i x_i. \quad (3)$$

Here, W denotes the mutual interactions of the spins, \mathbf{h} the outer field, and $e^{-\psi_q}$ is the normalization constant, $Z = e^{\psi_q}$ is the partition function, and $\psi_q = \psi_q(W, \mathbf{h})$ is called the free energy in physics or the cumulant generating function in statistics.

Let \mathcal{S} be the family of probability distributions of the above form (1), where (W, \mathbf{h}) forms a coordinate system to specify each distribution in \mathcal{S} . Let E_q be the expectation operator with respect to q . Then, the expectations of X and \mathbf{x} ,

$$K_q = E_q[X] = (E_q[x_i x_j]), \quad \mathbf{m}[q] = E_q[\mathbf{x}] = (E_q[x_i]). \quad (4)$$

form another coordinate system of \mathcal{S} . Our theory is based on information geometry and is applicable to many other general cases, but for simplicity's sake, we stick on this simple problem of obtaining a good approximation of $\mathbf{m}[q]$ for a given q .

Let us consider the subspace \mathbf{E} of \mathbf{S} such that x_i 's are independent or $W = 0$. A distribution $p \in \mathbf{E}$ is written as

$$p(x, \mathbf{h}) = \exp \{ \mathbf{h} \cdot \mathbf{x} - \psi_p \}. \quad (5)$$

This is a submanifold of \mathbf{S} specified by $W = 0$, and \mathbf{h} is its coordinates. The expectation

$$\mathbf{m} = \overline{E_p[\mathbf{x}]} \quad (6)$$

is another coordinate system of \mathbf{E} .

Physicists know that it is computationally difficult to calculate $\mathbf{m}[q]$ from $q(\mathbf{x}, W, \mathbf{h})$. It is given by

$$\mathbf{m}[q] = \frac{\partial}{\partial \mathbf{h}} \psi_q(W, \mathbf{h}) \quad (7)$$

but the partition function $Z_q = e^{-\psi_q}$ is difficult to obtain when the system size is large. On the other hand, for $p \in \mathbf{E}$, it is easy to obtain $\mathbf{m} = \overline{E_p[\mathbf{x}]}$ because x_i are independent. Hence, the mean field approximation tries to use quantities obtained in the form of expectation with respect to some relevant $p \in \mathbf{E}$.

Physicists established the method of approximation, called the mean field theory, including the TAP approximation. The problem is formulated more naturally in the framework of information geometry (Tanaka [this volume], Kappen [this volume] and Bhattacharyya and Keerthi [2000]). The present paper tries to give another way to elucidate this problem by information geometry of the α -connections introduced by Nagaoka and Amari [1982], Amari [1985] and Amari and Nagaoka [2000].

3 Concepts from Information Geometry

Here, we introduce some concepts of information geometry without entering in details. Let y be a discrete random variable taking values on a finite set $\{0, 1, \dots, N-1\}$, and let $p(y)$ and $q(y)$ be two probability distributions. In the case of spins, y represents $2^n \mathbf{x}$'s where $N = 2^n$.

α -divergence

The α -divergence from q to p is defined by

$$D_\alpha[q : p] = \frac{4}{1 - \alpha^2} \left(1 - \sum_y q^{\frac{1-\alpha}{2}} p^{\frac{1+\alpha}{2}} \right), \quad \alpha \neq \pm 1 \quad (8)$$

and for $\alpha = \pm 1$

$$D_{-1}[q : p] = \sum q \log \frac{q}{p}, \quad (9)$$

$$D_1[q : p] = \sum p \log \frac{p}{q}. \quad (10)$$

The latter two are the Kullback-Leibler divergence and its reverse. When $\alpha = 0$, it is the Hellinger distance,

$$D_0[q : p] = 2 \sum (\sqrt{p} - \sqrt{q})^2. \quad (11)$$

The divergence satisfies

$$D_\alpha[q : p] \geq 0, \quad (12)$$

with equality when and only when $q = p$. However, it is not symmetric except for $\alpha = 0$, and it satisfies

$$D_\alpha[q : p] = D_{-\alpha}[p : q]. \quad (13)$$

The α -divergence may be calculated in the following way. Let us consider a curve of probability distributions parameterized by t ,

$$p(y, t) = e^{-\psi(t)} q^{\frac{1-t}{2}} p^{\frac{1+t}{2}} = \exp \left\{ \frac{1}{2} \log pq + \frac{t}{2} \log \frac{p}{q} - \psi(t) \right\}, \quad (14)$$

which is an exponential family connecting q and p . Here, $e^{-\psi(t)}$ is the normalization constant. We then have

$$D_\alpha[q : p] = \frac{4}{1 - \alpha^2} (1 - e^{\psi(\alpha)}), \quad \alpha \neq \pm 1. \quad (15)$$

We also have

$$\psi'(\alpha) = \frac{1}{2} E_\alpha \left[\log \frac{p}{q} \right], \quad (16)$$

$$\psi''(\alpha) = E_\alpha \left[\frac{1}{4} \left(\log \frac{p}{q} \right)^2 \right] - \{\psi'(\alpha)\}^2, \quad (17)$$

$$\psi'''(\alpha) = \frac{1}{8} E_\alpha \left[\left\{ \left(\log \frac{p}{q} \right) - E_\alpha \left[\log \frac{p}{q} \right] \right\}^3 \right], \quad (18)$$

where E_α denotes expectation with respect to $p(y, \alpha)$. For $\alpha = \pm 1$, $\psi(\pm 1) = 0$. By taking the limit, we have

$$D_\alpha = 2\alpha\psi'(\alpha), \quad \alpha = \pm 1. \quad (19)$$

The family of the α -divergences gives invariant measures provided by information geometry.

4 Parametric model and Fisher information

When $p(y)$ is specified by a coordinate system ξ , it is written as $p(y, \xi)$. The $N - 1$ quantities

$$p_i = \text{Prob}\{y = i\}, \quad i = 1, \dots, N - 1 \quad (20)$$

form coordinates of $p(y)$. There are many other coordinate systems. For example

$$\theta_i = \log \frac{p_i}{p_0} \quad (21)$$

is another coordinate system. Let

$$\delta_i(y) = \begin{cases} 1, & y = i, \\ 0, & \text{otherwise.} \end{cases} \quad (22)$$

Then, we have

$$p(y, \xi) = \sum p_i \delta_i(y) + p_0 \delta_0(y), \quad (23)$$

$$p(y, \theta) = \exp \left\{ \sum \theta_i \delta_i(y) + \log p_0 \right\}. \quad (24)$$

The α -divergence for two nearby distributions $p(y, \xi)$ and $p(y, \xi + d\xi)$ is expanded as

$$D_\alpha [p(y, \xi) : p(y, \xi + d\xi)] = \frac{1}{2} \sum_{i,j} g_{ij}(\xi) d\xi_i d\xi_j, \quad (25)$$

where the right-hand side does not depend on α . The matrix $G(\xi) = (g_{ij}(\xi))$ is positive-definite and symmetric, given by

$$g_{ij}(\xi) = E \left[\frac{\partial \log p(y, \xi)}{\partial \xi_i} \frac{\partial \log p(y, \xi)}{\partial \xi_j} \right]. \quad (26)$$

This is called the Fisher information. It gives the unique invariant Riemannian metric to the manifold of probability distributions.

α -projection

Let \mathbf{M} be a submanifold in \mathbf{S} . Given $q \in \mathbf{S}$, the point $p^* \in \mathbf{M}$ is called the α -projection of q to \mathbf{M} , when function $D_\alpha[q, p], p \in \mathbf{M}$ takes a critical value at p^* , that is

$$\frac{\partial}{\partial \xi} D_\alpha[q : p(y, \xi)] = 0, \quad (27)$$

at p^* where ξ is a coordinate system of \mathbf{M} . The minimizer of $D_\alpha[q : p], p \in \mathbf{M}$, is the α -projection of q to \mathbf{M} . We denote it by

$$p^* = \prod_{\alpha} q. \quad (28)$$

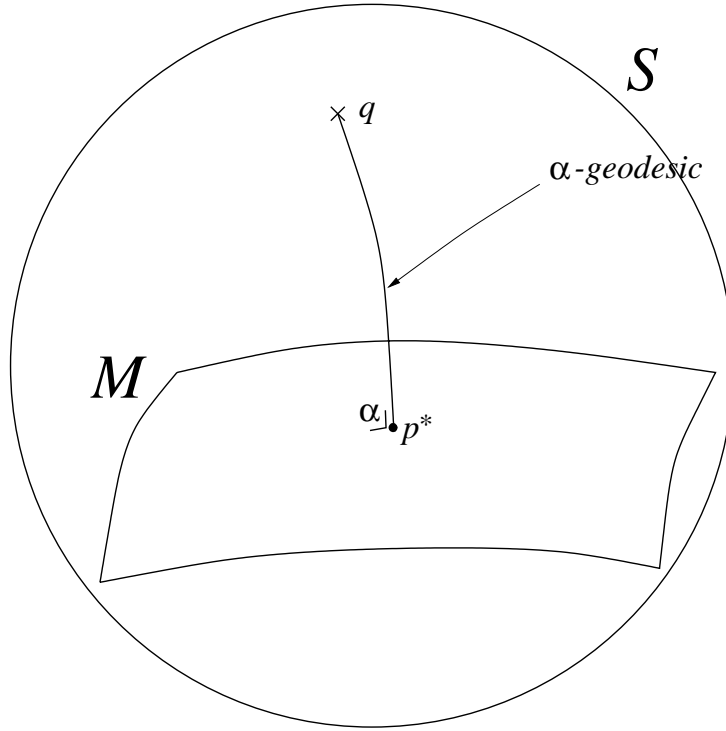


Figure 1: α -projection

In order to characterize the α -projection, we need to define the α -affine connection and α -geodesic derived therefrom. We do not explain them here (see Amari, 1985; Amari and Nagaoka, 2000). We show the following fact. See Fig.1.

Theorem 1. A point $p^* \in \mathbf{M}$ is the α -projection of q to \mathbf{M} , when and only when the α -geodesic connecting q and p^* is orthogonal to \mathbf{M} in the sense of the Fisher Riemannian metric G .

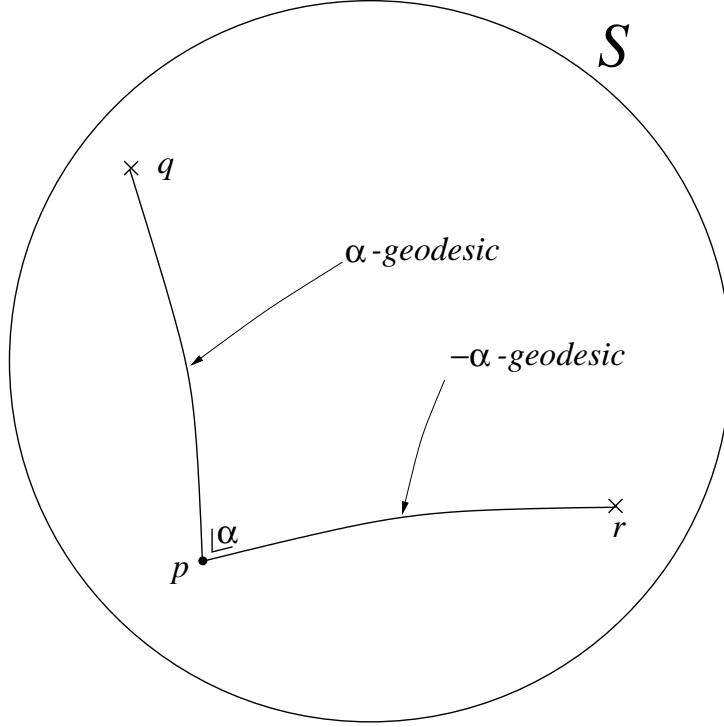


Figure 2: Pythagoras' theorem

Exponential family

A family of distributions is called an exponential family, when its probability distributions are written as

$$p(y, \theta) = \exp \left\{ \sum \theta_i k_i(y) - \psi(\theta) \right\} \quad (29)$$

by using an appropriate coordinate system θ , where $\mathbf{k} = k_i(y)$ are adequate functions of y . The spin system or Boltzmann machine (1) is an exponential family, where

$$\theta = (W, \mathbf{h}) \quad (30)$$

and \mathbf{k} consists of

$$\mathbf{k} = (X, \mathbf{x}). \quad (31)$$

The exponential family forms an $\alpha = \pm 1$ flat manifold, that is, $\alpha = \pm 1$ Riemann-Christoffel curvatures vanish identically, but this is a non-Euclidean space. There exist $\alpha = \pm 1$ affine coordinate systems in such a manifold. The above θ is an $\alpha = 1$ affine coordinate system, called the e-affine (exponential-affine), because the log probability is

linear in θ . An e-geodesic ($\alpha = 1$ geodesic) is linear in θ . More generally, for any two distributions $p(y)$ and $q(y)$, the e-geodesic connecting them is given by

$$\log p(y, t) = (1 - t) \log p(y) + t \log q(y) - \psi(t). \quad (32)$$

Let us denote the expectation of \mathbf{k} by η ,

$$\eta = E[\mathbf{k}]. \quad (33)$$

It is known that this η forms another coordinate system of an exponential family. This is an $\alpha = -1$ affine coordinate system, or m-affine (mixture affine) coordinate system. The two coordinate systems are connected by the Legendre transformation,

$$\eta = \frac{\partial}{\partial \theta} \psi(\theta), \quad (34)$$

$$\theta = \frac{\partial}{\partial \eta} \varphi(\eta), \quad (35)$$

where $\varphi(\eta)$ is the negative of entropy function, and

$$\psi(\theta) + \varphi(\eta) - \theta \cdot \eta = 0 \quad (36)$$

holds. Any linear curve in η is an m-geodesic.

An important property is given by the following theorem. See Fig.2.

Theorem 2. Let p, q, r be three probability distributions in an $\pm\alpha$ -flat manifold \mathbf{S} . When the α -geodesic connecting p and q is orthogonal at q with respect to the Riemannian metric to the $-\alpha$ -geodesic connecting q and r ,

$$D_\alpha[p : q] + D_\alpha[q : r] = D_\alpha[p : r]. \quad (37)$$

From this follows

Theorem 3. Let \mathbf{M} be a smooth submanifold in an $\pm\alpha$ -flat manifold \mathbf{S} , and let p^* be the α -projection from q to \mathbf{M} . Then, the α -geodesic connecting q and p^* is orthogonal to \mathbf{M} and vice versa.

5 Geometry of \mathbf{E}

Since \mathbf{E} consists of all the independent distributions, it is easy to show the geometry of \mathbf{E} . Moreover, \mathbf{E} itself is an exponential family,

$$p(\mathbf{x}, \bar{\mathbf{h}}) = \exp \{ \bar{\mathbf{h}} \cdot \mathbf{x} - \psi_0(\bar{\mathbf{h}}) \}, \quad (38)$$

where

$$e^{\psi_0(\bar{\mathbf{h}})} = \prod_i (e^{\bar{h}_i} + e^{-\bar{h}_i}) \quad (39)$$

or

$$\psi_0(\bar{\mathbf{h}}) = \sum_i \log (e^{\bar{h}_i} + e^{-\bar{h}_i}). \quad (40)$$

This $\bar{\mathbf{h}} = (\bar{h}_1, \dots, \bar{h}_n)$ is the e-affine coordinates of \mathbf{E} . Its m-affine coordinates are given by

$$\mathbf{m} = E_p[\mathbf{x}] = \frac{\partial}{\partial \bar{\mathbf{h}}} \psi_0(\bar{\mathbf{h}}), \quad (41)$$

which is easily calculated as

$$m_i = E_p[x_i] = \frac{\partial}{\partial \bar{h}_i} \psi_0(\bar{\mathbf{h}}) = \frac{e^{\bar{h}_i} - e^{-\bar{h}_i}}{e^{\bar{h}_i} + e^{-\bar{h}_i}} = \tanh \bar{h}_i. \quad (42)$$

This is solved as

$$e^{\bar{h}_i} = \sqrt{\frac{1 + m_i}{1 - m_i}}. \quad (43)$$

In terms of \mathbf{m} , the probability is written as

$$p(\mathbf{x}, \mathbf{m}) = \prod \frac{1 + m_i x_i}{2}, \quad x_i = \pm 1. \quad (44)$$

The Riemannian metric or the Fisher information $G = (g_{ij})$ is

$$\begin{aligned} g_{ij} &= \frac{\partial m_i}{\partial \bar{h}_j} \\ &= (1 - m_i^2) \delta_{ij}. \end{aligned} \quad (45)$$

Its inverse $\bar{G} = G^{-1} = (\bar{g}_{ij})$ is

$$\bar{g}_{ij} = \left(\frac{1}{1 - m_i^2} \right) \delta_{ij}. \quad (46)$$

Let $l(\mathbf{x}, \mathbf{m}) = \log p(\mathbf{x}, \mathbf{m})$. We then have

$$\partial_{m_i} l = \frac{x_i - m_i}{1 - m_i^2}, \quad (47)$$

$$\partial_{m_i}^2 l = -(\partial_{m_i} l)^2. \quad (48)$$

6 The α -projection and mean field approximation

Given $q \in S$, its α -projection to \mathbf{E} is given by

$$\bar{p}_\alpha = \prod_{\alpha} q = \arg \min_{p \in \mathbf{E}} D_\alpha [q : p]. \quad (49)$$

We denote by $\mathbf{m}_\alpha[q]$ the expectation of \mathbf{x} with respect to \bar{p}_α , that is $E_{\bar{p}_\alpha}[x]$. Then it is given by

$$\frac{\partial}{\partial \mathbf{m}} D_\alpha [q : p(\mathbf{x}, \mathbf{m}_\alpha)] = 0. \quad (50)$$

From the point of view of information geometry, $\prod_{\alpha} q = p(\mathbf{x}, \mathbf{m}_\alpha) \in \mathbf{E}$ is the α -geodesic projection of q to \mathbf{E} in the sense that the α -geodesic connecting q and p is orthogonal to \mathbf{E} at $p = p(\mathbf{x}, \mathbf{m}_\alpha)$.

When $\alpha = -1$, \bar{p}_{-1} is the m -projection ($\alpha = -1$ -projection) of q to \mathbf{E} . We have

$$\mathbf{m}_{-1} = \mathbf{m}[q] \quad (51)$$

which is the quantity we want to obtain. This relation is directly calculated by solving

$$\frac{\partial}{\partial \mathbf{m}} D_{-1} [q : p(\mathbf{m}_{-1})] = 0, \quad (52)$$

because this is equivalent to

$$\frac{\partial}{\partial \mathbf{m}} \int q \log p(\mathbf{x}, \mathbf{m}) dx = 0, \quad (53)$$

or

$$\frac{\partial}{\partial \mathbf{m}} E_q [l(\mathbf{x}, \mathbf{x})] = 0 = \text{const} E_q [\mathbf{x} - \mathbf{m}]. \quad (54)$$

Hence,

$$\mathbf{m}_{-1} = E_q[\mathbf{x}] \quad (55)$$

which is the quantity we have searched for. But we cannot calculate $E_q[\mathbf{x}]$ explicitly, due to the difficulty in calculating Z or ψ_q for q .

Physicists tried to obtain \mathbf{m}_{-1} by mean field approximation in an intuitive way. If we use the e-projection of q to \mathbf{E} instead of the m -projection, we have the naive mean field

approximation (Tanaka, 2000). To show this, we calculate the e-projection (1-projection) of q to \mathbf{E} . For $\alpha = 1$,

$$\begin{aligned} D_1[q : p] &= D_{-1}[p : q] = E_p [\bar{\mathbf{h}} \cdot \mathbf{x} - \psi_p - (WX + \mathbf{h} \cdot \mathbf{x} - \psi_q)] \\ &= \bar{\mathbf{h}} \cdot \mathbf{m} - \psi_p - W \cdot M - \mathbf{h} \cdot \mathbf{m} + \psi_q \end{aligned} \quad (56)$$

because of $M = E_p[X] = \mathbf{m}\mathbf{m}^T$. Hence,

$$\begin{aligned} \frac{\partial D_1}{\partial \mathbf{m}} &= \frac{\partial \bar{\mathbf{h}}}{\partial \mathbf{m}} \frac{\partial}{\partial \bar{\mathbf{h}}} (\bar{\mathbf{h}} \cdot \mathbf{m} - \psi_p) - W\mathbf{m} - \mathbf{h} \\ &= \tanh^{-1}\mathbf{m} - W\mathbf{m} - \mathbf{h}. \end{aligned} \quad (57)$$

This gives

$$\mathbf{m}_1[q] = \tanh [W\mathbf{m}_1[q] + \mathbf{h}], \quad (58)$$

known as the “naive” mean-field approximation. In the component form, this is

$$m_i = \tanh \left(\sum w_{ij} m_j + h_i \right). \quad (59)$$

This equation can have a number of solutions. It is necessary to check which solution minimizes $D_1[q : p]$. The minimization may be attained at the boundary of $m_i = \pm 1$.

Similarly, we have the α -projection $\mathbf{m}_\alpha[q]$ by solving

$$\frac{\partial}{\partial \mathbf{m}} D_\alpha [q : p(\mathbf{x}, \mathbf{m})] = 0. \quad (60)$$

However, it is not tractable to obtain \mathbf{m}_α explicitly except for $\alpha = 1$.

7 α -trajectory

For a fixed q , its α -projection $\mathbf{m}_\alpha[q]$ is considered as a path in \mathbf{E} connecting the true $\mathbf{m}_{-1} = \mathbf{m}[q]$ and the mean-field approximation $\mathbf{m}_1[q]$. This is called the α -trajectory of q in \mathbf{E} . See Fig.3.

The tangent direction of the trajectory is given by $d\mathbf{m}_\alpha/d\alpha$. This is given from

$$\partial_m D_\alpha [q : p(\mathbf{m}_\alpha)] = 0, \quad (61)$$

$$\partial_m D_{\alpha+d\alpha} [q : p(\mathbf{m}_\alpha + d\mathbf{m}_\alpha)] = 0 \quad (62)$$

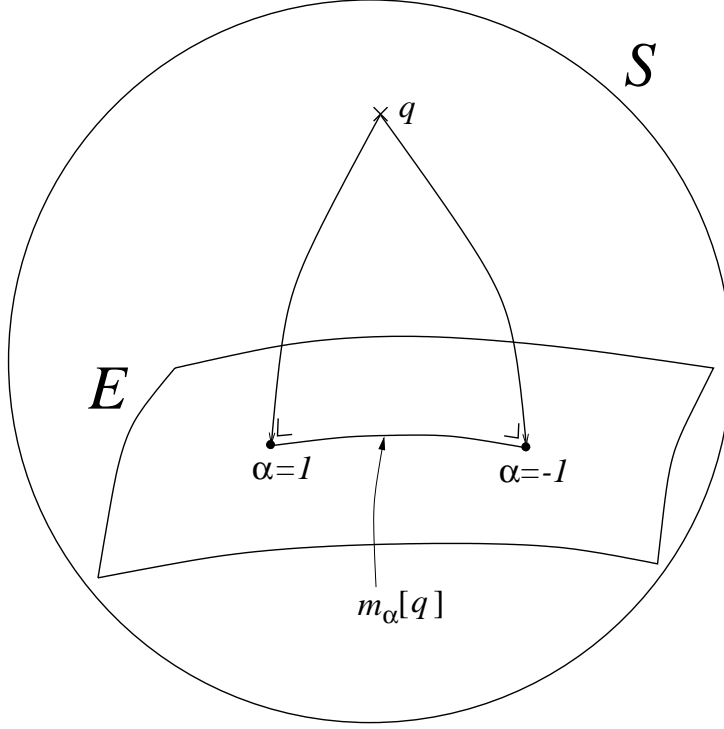


Figure 3: α -trajectory

so that, by Taylor expansion,

$$\partial_m \partial_\alpha D_\alpha d\alpha + \partial_m^2 D_\alpha d\mathbf{m}_\alpha = 0. \quad (63)$$

Here, $\partial_m = \partial/\partial \mathbf{m}$ and $\partial_\alpha = d/d\alpha$. We then have

$$\frac{d\mathbf{m}_\alpha}{d\alpha} = - \{ \partial_m^2 D_\alpha \}^{-1} \partial_m \partial_\alpha D_\alpha. \quad (64)$$

Starting from the naive approximation, we may improve it by the expansion

$$\mathbf{m}_\alpha[q] = \mathbf{m}_1[q] + \frac{d\mathbf{m}_\alpha}{d\alpha}(\alpha - 1) + \frac{1}{2} \frac{d^2 \mathbf{m}_\alpha}{d\alpha^2}(\alpha - 1)^2 + \dots \quad (65)$$

where the derivatives are evaluated at $\alpha = 1$, provided they are calculated easily. Another idea is to integrate $d\mathbf{m}_\alpha/d\alpha$, provided the derivative at α can be calculated. We cannot solve these methods now.

In order to study the α -trajectory, we show some preliminary calculations. The α -divergence from q to $p(x, \mathbf{m}) \in \mathbf{E}$ is written as

$$D_\alpha [q : p(\mathbf{m})] = \frac{4}{1 - \alpha^2} (1 - e^{\psi(\alpha, \mathbf{m})}) \quad (66)$$

where

$$\psi(\alpha, \mathbf{m}) = \log \sum q^{\frac{1-\alpha}{2}} p^{\frac{1+\alpha}{2}} = \log E_p \left[\left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right]. \quad (67)$$

We first calculate $\partial_m \psi$, since \mathbf{m}_α is given by

$$\partial_m \psi(\alpha, \mathbf{m}_\alpha) = 0. \quad (68)$$

We have

$$\begin{aligned} \partial_m \psi &= \frac{1+\alpha}{2} e^{-\psi(\alpha, \mathbf{m})} \sum q^{\frac{1-\alpha}{2}} p^{-\frac{1-\alpha}{2}} \partial_m p \\ &= \frac{1+\alpha}{2} e^{-\psi} E_p \left[\partial_m l \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right]. \end{aligned} \quad (69)$$

We then have

$$\partial_m^2 \psi = -(\partial_m \psi)^2 + \frac{1+\alpha}{2} e^{-\psi} \partial_m E_p \left[\partial_m l \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right], \quad (70)$$

where $\partial_m^2 \psi$ is a matrix and $(\partial_m \psi)^2$ implies $(\partial_m \psi)(\partial_m \psi)^T$. At $\mathbf{m} = \mathbf{m}_\alpha$, the first term vanishes and

$$\partial_m^2 \psi = \frac{1+\alpha}{2} e^{-\psi} E_p \left[((\partial_m l)^2 + \partial_m^2 l) \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} - \frac{1-\alpha}{2} (\partial_m l)^2 \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right]. \quad (71)$$

We also have

$$\begin{aligned} \partial_m \partial_\alpha \psi(\alpha, \mathbf{m}) &= \frac{1+\alpha}{4} e^{-\psi} E_p \left[\partial_m l \log \frac{p}{q} \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right] + \frac{1}{1+\alpha} \partial_m \psi \\ &= \frac{1+\alpha}{4} e^{-\psi} E_p \left[\partial_m l \log \frac{p}{q} \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right]. \end{aligned} \quad (72)$$

From this we have

$$\frac{d\mathbf{m}_\alpha}{d\alpha} = -\frac{1}{2} A^{-1} E_p \left[\partial_m l \log \frac{p}{q} \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right] \quad (73)$$

$$A = E_p \left[((\partial_m l)^2 + \partial_m^2 l) \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} - \frac{1-\alpha}{2} (\partial_m l)^2 \left(\frac{q}{p} \right)^{\frac{1-\alpha}{2}} \right]. \quad (74)$$

For $\alpha = 1$, we have

$$\frac{d\mathbf{m}_\alpha}{d\alpha} = -\frac{1}{2} \frac{E_p [\partial_m l \{\log(p/q)\}^2]}{E_p [\{(\partial_m l)^2 + \partial_m^2 l\} \log(p/q) + (\partial_m l)^2]}. \quad (75)$$

8 α -Hessian

The α -projection $\mathbf{m}_\alpha[q]$ is given by the point \bar{p}_α in \mathbf{E} that is the orthogonal projection of q to \mathbf{E} by the α -geodesic. Such projection is not necessarily unique. The projection is not necessarily the minimizer of $D_\alpha[q : p]$ but is a saddle or even the maximizer of D_α . To elucidate this we calculate the Hessian $H^\alpha = (H_{ij}^\alpha)$

$$H_{ij}^\alpha = \frac{\partial^2}{\partial m_i \partial m_j} D_\alpha[q : p(\mathbf{m})] \quad (76)$$

at the α -projection $\mathbf{m}_\alpha[q]$. When H^α is positive-definite, the α -projection \mathbf{m}_α gives a local minimum, but it is otherwise a saddle or local maximum. We have

$$\begin{aligned} H_{ij}^\alpha &= -\frac{2}{1+\alpha} E_p \left[\left\{ \frac{1+\alpha}{2} \partial_i l \partial_j l + \partial_i \partial_j l \right\} f_\alpha \right] \\ &= -E_p \left[\left\{ \partial_i l \partial_j l + \frac{2}{1+\alpha} \partial_i \partial_j l \right\} f_\alpha \right], \end{aligned} \quad (77)$$

where $\partial_i = \partial/\partial m_i$ and $f_\alpha = \frac{2}{1-\alpha} \left(\frac{q}{p}\right)^{\frac{1-\alpha}{2}}$. From

$$\partial_i l = \frac{1}{1-m_i^2} (x_i - m_i), \quad (78)$$

$$\partial_i \partial_j l = -\delta_{ij} \frac{1}{(1-m_i^2)^2} (x_i - m_i)^2, \quad (79)$$

we finally have

$$\begin{aligned} H_{ij}^\alpha &= \frac{-1}{(1-m_i^2)(1-m_j^2)} E_p [(x_i - m_i)(x_j - m_j) f_\alpha] \\ &= -\bar{g}_{ii} \bar{g}_{jj} \{E_p [x_i x_j f_\alpha] - m_i m_j\}, \quad i \neq j \end{aligned} \quad (80)$$

because of

$$m_i^\alpha[q] = E_p [x_i f_\alpha] \quad (81)$$

and for $i \neq j$

$$H_{ii}^\alpha = \frac{1-\alpha}{1+\alpha} (\bar{g}_{ii})^2 E_p [(x_i - m_i)^2 f_\alpha]. \quad (82)$$

We calculate the two special cases $\alpha = \pm 1$. For $\alpha = -1$,

$$\begin{aligned} H_{ij}^{-1} &= \partial_i \partial_j \int q \log \frac{q}{p} dx \\ &= (\bar{g}_{ii})^2 \delta_{ij} E_q [(x_i - m_i)^2] \\ &= \bar{g}_{ii} \delta_{ij} = \bar{G}, \end{aligned} \quad (83)$$

where \bar{G} is the inverse of the Fisher information matrix of \mathbf{E} . This is diagonal and positive-definite. Because \mathbf{E} is e-flat, we know that $\alpha = -1$ -projection gives the global minimum, is unique (that is no other critical points) and gives the true solution $\mathbf{m}_{-1}[q] = E_q[\mathbf{x}]$.

For $\alpha = 1$, we have

$$\begin{aligned} H_{ij}^1 &= \partial_i \partial_j \int p \log \frac{p}{q} dx \\ &= \bar{g}_{ii} \delta_{ij} - E_p [(\partial_i \partial_j l + \partial_i l \partial_j l) \log q]. \end{aligned} \quad (84)$$

Hence,

$$\begin{aligned} H_{ii}^1 &= \bar{g}_{ii} = \frac{1}{1 - m_i^2}, \\ H_{ij}^1 &= -\bar{g}_{ii} \bar{g}_{jj} E_p [(x_i - m_j)(x_j - m_j) \log q] \\ &= -\bar{g}_{ii} \bar{g}_{jj} E_p \left[(x_i - m_i)(x_j - m_j) \left\{ \sum w_{kl} x_k x_l \right. \right. \\ &\quad \left. \left. + \sum h_k x_k - \psi_q \right\} \right] \\ &= -w_{ij}. \end{aligned} \quad (85)$$

This shows that H^1 is not necessarily positive-definite. This fact is related to the α -curvature of \mathbf{E} . When $n = 2$ (two neurons), it is positive definite when and only when $w = w_{12}$ satisfies

$$w^2 < \frac{1}{1 - m_1^2} \frac{1}{1 - m_2^2}. \quad (87)$$

Otherwise, it is a saddle. When $h_1 = h_2$, there exist one or two local minima other than this.

This fact implies that the naive mean field approximation might give a pathological solution in some cases. When $|w|$ is large, the above two-spin system is dynamically bistable, having two stable solutions $x_1 = x_2 = 1$, $x_1 = x_2 = -1$ ($w > 0$) or $x_1 = -x_2$ ($w < 0$).

When $\alpha = -1$, the α -projection \mathbf{m}_{-1} is unique. Starting from $\alpha = -1$, the α -trajectory \mathbf{m}_α bifurcates at some α , and then bifurcates further, depending on W . See Fig.4. When W is large, the naive mean field approximation (59) may have an exponentially large number of solutions. It is interesting to study the diagram of bifurcation for the α -trajectory.

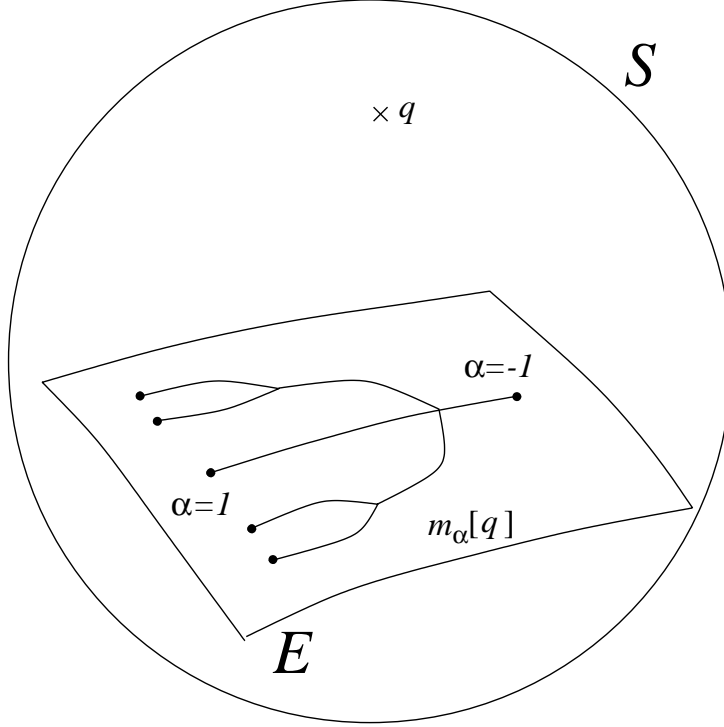


Figure 4: Bifurcation

9 Small w approximation of the α -trajectory

We give an explicit formula for the α -projection \mathbf{m}_α assuming that w_{ij} are small. The TAP solution corresponds to $\mathbf{m}_{-1}[q]$ under this approximation.

Let us consider an exponential family $\{p(\mathbf{x}, \theta)\}$. For two nearby distributions $q = p(\mathbf{x}, \theta + d\theta)$ and $p = p(\mathbf{x}, \theta)$, we have

$$D_\alpha(q, p) = D_\alpha(\theta + d\theta, \theta) = \frac{1}{2} \sum g_{ij} d\theta^i d\theta^j + \frac{3-\alpha}{12} \sum T_{ijk} d\theta^i d\theta^j d\theta^k, \quad (88)$$

where

$$g_{ij} = \frac{\partial^2}{\partial \theta^i \partial \theta^j} \psi(\theta), \quad T_{ijk} = \frac{\partial^3}{\partial \theta^i \partial \theta^j \partial \theta^k} \psi(\theta) = \frac{\partial}{\partial \theta^k} g_{ij}. \quad (89)$$

In our case, $\theta = (W, \mathbf{h})$, and for $q = q(\mathbf{x}, dW, \mathbf{h})$ and $p = p(\mathbf{x}, 0, \bar{\mathbf{h}})$, we have

$$d\theta = (dW, d\mathbf{h}), \quad (90)$$

where $W = dW$ is assumed to be small and $d\mathbf{h} = \mathbf{h} - \bar{\mathbf{h}}$. We can calculate g_{ij} and T_{ijk} at $p \in \mathbf{E}$, for example, for $I = (i, j)$ and $J = (k, l)$,

$$g_{IJ} = E_p [(x_i x_j - m_i m_j) (x_k x_l - m_k m_l)]. \quad (91)$$

The metric G consists of the three parts g_{IJ} , g_{Ik} , g_{kl} where I, J etc are index pairs corresponding to $dW^I = w^{ij}$, and small letter indices i, j etc refer to $dh^i = h^i - \theta^i$.

Note that

$$g_{ij} = E[(x_i - m_i)(x_j - m_j)] \quad (92)$$

$$T_{ijk} = E[(x_i - m_i)(x_j - m_j)(x_k - m_k)] \quad (93)$$

etc. By using this, we have the following expansion,

$$\begin{aligned} D_\alpha[q, p] = & \frac{1}{2}g_{IJ}dw^I dw^J + \frac{1}{2}g_{ij}dh^i dh^j + g_{Ik}dw^I dh^k \\ & + \frac{3-\alpha}{12} \left\{ T_{IJK}dw^I dw^J dw^K + 3T_{Ijk}dw^I dh^j dh^k \right. \\ & \left. + 3T_{IJK}dw^I dw^J dh^k + T_{ijk}dh^i dh^j dh^k \right\}, \end{aligned} \quad (94)$$

where the summation convention is used for repeated indices. In order to obtain the α -projection, we solve

$$\frac{\partial D_\alpha}{\partial \bar{h}_i} = 0, \quad (95)$$

where indices i of $d\theta_i$ are decomposed into indices pairs $I = (i, j)$, etc. for $W^I = w^{ij}$ and single indices i, j, \dots for h_i . For example, we have

$$\begin{aligned} 0 = & -g_{il}dh^i - g_{Ii}dw^I - \frac{1-\alpha}{4} \left(T_{IJI}dw^I dw^J \right. \\ & \left. + T_{ijl}dh^i dh^j + 2T_{Ikl}dw^I dh^k \right) + O(w^3), \end{aligned} \quad (96)$$

where we used

$$\frac{\partial}{\partial \bar{h}_l} dh^i = -\delta_{il}. \quad (97)$$

The first-order solution to (96) does not depend on α ,

$$h^i - \bar{h}_i = g^{il}g_{Ii}dw^I \quad (98)$$

or

$$\theta^i = h^i + \sum w_{ij}m_j \quad (99)$$

which is the naive mean field approximation.

In order to calculate higher-order corrections, we note

$$g_{ij} = \delta_{ij} (1 - m_i^2) \quad (100)$$

$$g_{Ik} = (1 - m_i^2) \{ \delta_{ki} m_j + \delta_{kj} m_i \}, \quad I = (j, j) \quad (101)$$

$$g_{IJ} = E [(x_i x_j - m_i m_j) (x_k x_l - m_k m_l)], \quad I = (i, j), J = (k, l) \quad (102)$$

etc. Quantities T can be calculated similarly. The second-order correction term A_l , which is given by the second term of (96), is obtained after painful calculations as

$$A_l = m_l \sum (w_{lk})^2 (1 - m_k^2). \quad (103)$$

Some easy terms are

$$T_{ijl} dh^i dh^j = -2m_l (1 - m_l^2) (dw^{lk} m_k)^2 \quad (104)$$

$$\begin{aligned} T_{Ikl} dw^I dh^l &= 2m_l (1 - m_l^2) (dw^{lk} m_k)^2 \\ &- (1 - m_l^2) (1 - m_k^2) dw^{lk} dw^{ks} m_s. \end{aligned} \quad (105)$$

After all, we have

$$\bar{h}_l = m_l + \sum w_{lk} m_k + \frac{1 - \alpha}{2} m_l \sum (w_{lk})^2 (1 - m_k^2). \quad (106)$$

or

$$m_l^\alpha = \tanh \left(h_l + \sum w_{ek} m_k^\alpha + \frac{1 - \alpha}{2} m_l^\alpha \sum (w_{lk})^2 (1 - m_k^{\alpha 2}) \right). \quad (107)$$

This gives the α -projections in terms of parameter α , where $\alpha = 1$ is for the naive approximation and $\alpha = -1$ is for the TAP approximation. This is small w approximation of the α -trajectory and is valid for small w .

Conclusions

The present article studies the geometrical structure underlying mean field approximation. Information geometry is used for this purpose which has the Riemannian metric together with dual pairs of affine connections. Information geometry gives the α -structure to the manifold of probability distributions of the SK-spin glass system or the Boltzmann

machine. The α -divergence is defined in the manifold which is invariant under a certain criterion.

The mean field approximation is a method of calculating quantities related to a complex probability distribution, by using a simple tractable model such as the family of independent distributions. The $\alpha = -1$ projection of the distribution to submanifold consisting of independent distributions is known to give the correct answer, but it is intractable. The $\alpha = 1$ approximation is tractable, but it gives only an approximation.

We search for possibility of using the α -approximation that minimizes the α -divergence. It is unfortunately difficult to calculate. But its properties are studied for future study. We have also shown the information-geometric meaning of the TAP approximation.

We have elucidated the fact that $\alpha = -1$ projection is unique, giving the true solution but α -approximation ($\alpha \neq -1$) is not necessarily unique. When we study the trajectory consisting of the α -projections, there are a number of bifurcations where the α -projection bifurcates. It is an interesting problem to study the properties of such bifurcation and its implications.

The present article is preliminary to further studies on interesting problems connecting information geometry and statistical physics.

References

- [1] S. Amari, *Differential-Geometrical Methods in Statistics*, Lecture Notes in Statistics 28, Springer-Verlag, 1985.
- [2] S. Amari, Dualistic geometry of the manifold of higher-order neurons, *Neural Networks*, 4:443–451, 1991.
- [3] S. Amari, Information geometry of EM and em algorithms for neural networks, *Neural Networks*, 8:1379–1408, 1995.
- [4] S. Amari, Natural gradient works efficiently in learning, *Neural Computation*, 10:251–276, 1998.

- [5] S. Amari, K. Kurata, and H. Nagaoka, Information geometry of Boltzmann machines, *IEEE Transactions on Neural Networks*, 3:26—271, 1992.
- [6] S. Amari and H. Nagaoka, *Methods of Information Geometry*, AMS and Oxford University Press, 2000.
- [7] C. Bhattacharyya and S.S. Keerthi, Mean-field methods for Stochastic Connectionist Networks, TR No. Iisc-CSA-00-03.
- [8] C. Bhattacharyya and S.S. Keerthi, Plefka’s mean-field theory from a Variational viewpoint, TR NO. IISc-CSA-00-02.
- [9] H. Chernoff, A measure of asymptotic efficiency for tests of a hypothesis based on a sum of observations, *Annals of Mathematical Statistics*, 23:493–507, 1952.
- [10] N.N. Chentsov (Čencov), *Statistical Decision Rules and Optimal Inference*, American Mathematical Society, Rhode Island, U.S.A., 1982, (Originally published in Russian, Nauka, Moscow, 1972).
- [11] I. Csiszár, I -divergence geometry of probability distributions and minimization problems, *The Annals of Probability*, 3:146–158, 1975.
- [12] H. Nagaoka and S. Amari, Differential geometry of smooth families of probability distributions, Technical Report METR 82-7, Dept. of Math. Eng. and Instr. Phys, Univ. of Tokyo, 1982.
- [13] A. Rényi, On measures of entropy and information, In *Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, pages 547–561, University of California Press, 1961.
- [14] T. Tanaka, Information geometry of mean field approximation, *Neural Computation*, to appear.
- [15] C. Tsallis, Possible generalization of Boltzmann-Gibbs statistics, *Journal of Statistical Physics*, 52:479–, 1988.