

# Information Geometry for Turbo Decoding

Shiro Ikeda,<sup>1</sup> Toshiyuki Tanaka,<sup>2</sup> and Shun-ichi Amari<sup>3</sup>

<sup>1</sup>Institute of Statistical Mathematics, Minato, 106-8569 Japan

<sup>2</sup>Tokyo Metropolitan University, Hachioji, 192-0397 Japan

<sup>3</sup>RIKEN BSI, Wako, 351-0198 Japan

## SUMMARY

Turbo codes are known as a class of error-correcting codes which have high error-correcting performance with efficient decoding algorithm. Characteristics of the iterative decoding algorithm have been studied in detail through a variety of numerical experiments, but theoretical results are still insufficient. In this paper, this issue is addressed from the information geometrical viewpoint. As a result, a mathematical framework for analyzing turbo codes is obtained, and some of the fundamental properties of turbo decoding are elucidated based on this framework. Recently, it has been pointed out that the turbo decoding algorithm is related to the decoding algorithm of low-density parity check codes, the computation method of Bethe approximation in statistical physics, and the belief propagation algorithm of Bayesian networks. The mathematical framework given in the present paper can also be used to analyze these wide classes of iterative computation methods, and hence represent a new analysis tool. © 2004 Wiley Periodicals, Inc. *Syst Comp Jpn*, 36(1): 79–87, 2005; Published online in Wiley InterScience (www.interscience.wiley.com). DOI 10.1002/scj.10359

**Key words:** turbo code; MPM decoding; information geometry.

## 1. Introduction

Turbo codes are a class of error-correcting codes in which an iterative algorithm is used for decoding. Since their appearance in 1993 [1], a variety of numerical experiments have shown the high performance of the practical codes. However, theoretical results [2] reported are not enough, and we need further results to understand the fundamental properties of the codes.

Similarities between turbo decoding and other methods have also been noted. Common features were found [4] to exist between the iterative algorithm of the turbo decoding and the iterative algorithm used in the decoding of low-density parity check codes [3]. It was also shown that these decoding problems could be formularized as an inference problem of a Bayesian network, and that the iterative decoding algorithms were equivalent to belief propagation for this Bayesian network [5]. Common features were also shown to exist in the computation methods of Bethe approximation in statistical physics [6]. Of course, there are still many theoretically unresolved problems regarding these methods. Therefore, the framework of these iterative methods will become clear once the mathematical structure of turbo codes is clarified.

---

Contract grant sponsor: Supported in part by MEXT, Japan under a Grant-in-Aid for Scientific Research (16700227, 14084208).

In this paper, we first describe the turbo decoding algorithm from an information geometrical [7, 8] viewpoint, and build a mathematical framework to analyze it. The basic mathematical properties of the turbo decoding algorithm are also clarified based on this framework. The geometrical structure of turbo decoding results is elucidated, and local stability conditions of the algorithm are demonstrated. A cost function, which is important for turbo decoding, is also shown, and the decoding errors of turbo codes are discussed. The result of the present paper is also useful for the iterative algorithms of the previously described classes.

## 2. Information Geometry of Turbo Codes

### 2.1. Turbo codes

Let us consider a case in which an information data block  $\mathbf{x} = (x_1, \dots, x_N)^T$ ,  $x_i \in \{-1, +1\}$  is sent over a memoryless binary symmetric channel (BSC). In this paper, a BSC is assumed for simplicity, but the framework of the information geometry obtained in the present paper can be expanded to an additive white Gaussian noise channel and other memoryless channels. Turbo codes (Fig. 1) are implemented as convolutional codes in order to increase the code lengths, but we treat them as block codes in the present paper [2, 5]. Turbo codes use two encoders to create two parity check words for a single code word. The corresponding results are  $\mathbf{y}_1 = (y_{11}, \dots, y_{1L})^T$ ,  $\mathbf{y}_2 = (y_{21}, \dots, y_{2L})^T$ ,  $y_{1j}, y_{2j} \in \{-1, +1\}$ . When  $(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$  are transmitted over a communications channel, they are received as  $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ , where  $\tilde{x}_i, \tilde{y}_{1j}, \tilde{y}_{2j} \in \{-1, +1\}$ . Also,  $\mathbf{y}_r$ ,  $r = 1, 2$ , is a function of  $\mathbf{x}$ , and is expressed as  $\mathbf{y}_r(\mathbf{x})$  when needed. Based on the received words, the original information word is inferred as  $\hat{\mathbf{x}}$ .

A turbo-decoding algorithm will first be described. Turbo decoding is a type of decoding in which two decoders are used in an alternating manner. The probability distributions  $p(\tilde{\mathbf{x}}|\mathbf{x})$ ,  $p(\tilde{\mathbf{y}}_r|\mathbf{x})$ ,  $r = 1, 2$  and the variables shown below are defined, as is function  $F$ .

$$\begin{aligned} l_{x_i} &\stackrel{\text{def}}{=} \ln \frac{\sum_{\{\mathbf{x}:x_i=+1\}} p(\tilde{\mathbf{x}}|\mathbf{x})}{\sum_{\{\mathbf{x}:x_i=-1\}} p(\tilde{\mathbf{x}}|\mathbf{x})} \\ &= \ln \frac{p(\tilde{x}_i|x_i=+1)}{p(\tilde{x}_i|x_i=-1)}, \\ l_{y_{rj}} &\stackrel{\text{def}}{=} \ln \frac{\sum_{\{\mathbf{x}:y_{rj}=+1\}} p(\tilde{\mathbf{y}}_r|\mathbf{x})}{\sum_{\{\mathbf{x}:y_{rj}=-1\}} p(\tilde{\mathbf{y}}_r|\mathbf{x})} \\ &= \ln \frac{p(\tilde{y}_{rj}|y_{rj}=+1)}{p(\tilde{y}_{rj}|y_{rj}=-1)}, \end{aligned}$$

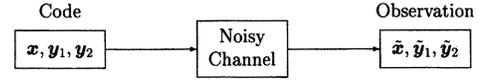


Fig. 1. Turbo codes.

$$\begin{aligned} L_r \mathbf{x} &\stackrel{\text{def}}{=} F(l_{\mathbf{x}}, l_{\mathbf{y}_r}) \\ &= \left\{ \ln \frac{\sum_{\{\mathbf{x}:x_i=+1\}} p(\tilde{\mathbf{x}}|\mathbf{x}) p(\tilde{\mathbf{y}}_r|\mathbf{x})}{\sum_{\{\mathbf{x}:x_i=-1\}} p(\tilde{\mathbf{x}}|\mathbf{x}) p(\tilde{\mathbf{y}}_r|\mathbf{x})} \right\} \end{aligned}$$

These are used to define the turbo decoding algorithm in the following manner (Fig. 2).

[Turbo Decoding]

1. Set  $\xi_1 = 0$ ,  $t = 1$ .
2. Calculate  $L_1 \mathbf{x}^{(t)} = F(l_{\mathbf{x}} + \xi_1, l_{\mathbf{y}_1})$ , and  $\xi_2$  is updated as

$$\xi_2 = L_1 \mathbf{x}^{(t)} - (l_{\mathbf{x}} + \xi_1) \quad (1)$$

3. Calculate  $L_2 \mathbf{x}^{(t)} = F(l_{\mathbf{x}} + \xi_2, l_{\mathbf{y}_2})$ , and  $\xi_1$  is updated as

$$\xi_1 = L_2 \mathbf{x}^{(t)} - (l_{\mathbf{x}} + \xi_2) \quad (2)$$

4. 2 and 3 are repeated while  $t$  is increased by 1 until  $L_1 \mathbf{x}^{(t)} = L_2 \mathbf{x}^{(t)} = L_1 \mathbf{x}^{(t+1)} = L_2 \mathbf{x}^{(t+1)}$  is satisfied.

The algorithm does not necessarily converge. The operations are repeated until it converges or a predetermined number of iterations is achieved. Usually this limit is from several to about ten.

### 2.2. MPM decoding

The goal of decoding turbo codes is to obtain the MPM (maximum posterior marginal) decoding. In MPM decoding, the distribution  $p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$  is marginalized down to each component of  $x_i$ , and the code which maximizes the product of marginal distributions is taken to be

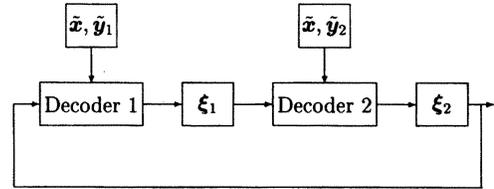


Fig. 2. Turbo decoding.

the inferred result. Let us first consider  $p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ . Since a memoryless BSC is assumed, the following is true:

$$p(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2|\mathbf{x}) = p(\tilde{\mathbf{x}}|\mathbf{x})p(\tilde{\mathbf{y}}_1|\mathbf{x})p(\tilde{\mathbf{y}}_2|\mathbf{x})$$

The right-hand side can be written as follows:

$$\begin{aligned} p(\tilde{\mathbf{x}}|\mathbf{x}) &= \exp(\beta\tilde{\mathbf{x}} \cdot \mathbf{x} - N\psi(\beta)) \\ p(\tilde{\mathbf{y}}_r|\mathbf{x}) &= \exp(\beta\tilde{\mathbf{y}}_r \cdot \mathbf{y}_r(\mathbf{x}) - L\psi(\beta)), \quad r = 1, 2 \\ \psi(\beta) &= \ln(e^{-\beta} + e^{\beta}) \end{aligned}$$

Here,  $\beta$  is a positive real number, and the bit error rate  $f_n$  of the BSC is expressed as  $f_n = (1 - \tanh\beta)/2$ . Assuming that  $c_0(\mathbf{x}) = \beta\tilde{\mathbf{x}} \cdot \mathbf{x}$ ,  $c_1(\mathbf{x}) = \beta\tilde{\mathbf{y}}_1 \cdot \tilde{\mathbf{y}}_1$ ,  $c_2(\mathbf{x}) = \beta\tilde{\mathbf{y}}_2 \cdot \tilde{\mathbf{y}}_2$ ,  $p(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2|\mathbf{x})$  is written as follows:

$$\begin{aligned} p(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2|\mathbf{x}) \\ = \exp(c_0(\mathbf{x}) + c_1(\mathbf{x}) + c_2(\mathbf{x}) - (N + 2L)\psi(\beta)) \end{aligned}$$

If we consider a uniform distribution  $p(\mathbf{x}) = 1/2^N$  as the prior distribution of  $\mathbf{x}$ , the posterior distribution of  $\mathbf{x}$  is given as follows:

$$\begin{aligned} p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) &= \frac{p(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2|\mathbf{x})}{\sum_{\mathbf{x}} p(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2|\mathbf{x})} \\ &= C \exp(c_0(\mathbf{x}) + c_1(\mathbf{x}) + c_2(\mathbf{x})) \quad (3) \end{aligned}$$

Also, we define

$$\Pi \circ p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) \stackrel{\text{def}}{=} \prod_{i=1}^N p_i(x_i|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$$

where  $\Pi$  is the operator of marginalization.

MPM decoding can be defined in the following manner:

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmax}} \Pi \circ p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$$

In turbo decoding, the MPM decoding is the ultimate goal, but the calculation cost of marginalization grows exponentially with respect to  $N$ , and the marginalization is not tractable. The turbo decoding algorithm provides an approximation to MPM decoding.

### 2.3. Preliminaries of information geometry

In this section, preliminaries of information geometry are given. Let us consider a set  $S$  of probability distributions for  $\mathbf{x}$ . It is a set of multinomial distributions of  $2^N$  elements. Its degree of freedom is  $(2^N - 1)$ , and it is an exponential family:

$$S = \left\{ p(\mathbf{x}) \mid p(\mathbf{x}) \geq 0, \mathbf{x} \in \{-1, +1\}^N, \sum p(\mathbf{x}) = 1 \right\}$$

Let us now define the  $e$ -flat and  $m$ -flat submanifolds included in  $S$ .

$e$ -flat: Manifold  $M \in S$  is  $e$ -flat when the  $r(\mathbf{x}; t)$  defined by the following equation is contained in  $M$ , where all  $q(\mathbf{x}), p(\mathbf{x}) \in M$ :

$$\ln r(\mathbf{x}; t) = (1 - t) \ln q(\mathbf{x}) + t \ln p(\mathbf{x}) + c, \quad t \in \mathcal{R}$$

Here,  $c$  is the normalization constant.

$m$ -flat: Manifold  $M \in S$  is  $m$ -flat when the  $r(\mathbf{x}; t)$  defined by the following equation is contained in  $M$ , where all  $q(\mathbf{x}), p(\mathbf{x}) \in M$ :

$$r(\mathbf{x}; t) = (1 - t)q(\mathbf{x}) + tp(\mathbf{x}), \quad t \in [0, 1]$$

An  $m$ -projection will next be defined.

[Definition 1] Let us assume that  $M$  is an  $e$ -flat submanifold of  $S$ . An  $m$ -projection from  $q(\mathbf{x}) \in S$  to  $M$  will be a point on  $M$  in which the *Kullback–Leibler (KL)* divergence from  $q(\mathbf{x})$  to  $M$  is at a minimum, and it can be defined in the following manner:

$$\Pi_M \circ q(\mathbf{x}) = \underset{p(\mathbf{x}) \in M}{\operatorname{argmin}} D[q(\mathbf{x}); p(\mathbf{x})]$$

[Theorem 1] The  $m$ -projection  $\Pi_M \circ q(\mathbf{x})$  of  $S$  from  $q(\mathbf{x}) \in S$  to an  $e$ -flat submanifold  $M$  is unique.

The KL divergence  $D[\cdot; \cdot]$  is defined as follows:

$$D[q(\mathbf{x}); p(\mathbf{x})] = \sum_{\mathbf{x}} q(\mathbf{x}) \ln \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

The KL divergence satisfies the inequality  $D[q(\mathbf{x}); p(\mathbf{x})] \geq 0$  and becomes 0 if  $q(\mathbf{x}) = p(\mathbf{x})$  holds for any  $\mathbf{x}$ .

To understand the turbo decoding algorithm, let us consider a submanifold  $M_D$  consisting of factorizable distribution. The definition will be as follows:

$$M_D = \left\{ p(\mathbf{x}; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta} \cdot \mathbf{x} - \psi(\boldsymbol{\theta})) \mid \boldsymbol{\theta} \in \mathcal{R}^N \right\}$$

where  $\psi(\boldsymbol{\theta})$  is the normalizing factor defined in the following manner:

$$\psi(\boldsymbol{\theta}) = \ln \sum_{\mathbf{x}} \exp(\boldsymbol{\theta} \cdot \mathbf{x}) = \sum_i \ln(e^{-\theta_i} + e^{\theta_i})$$

By its definition,  $M_D$  is an exponential family, and since an exponential family is  $e$ -flat,  $M_D$  is an  $e$ -flat submanifold [8]. Parameter  $\boldsymbol{\theta}$  gives the coordinate system of the manifold  $M_D$  and is called the natural parameter. Another coordinate system,  $\boldsymbol{\eta}$ , which is called the expectation parameter, is defined as follows:

$$\boldsymbol{\eta} = \sum_{\mathbf{x}} p(\mathbf{x}; \boldsymbol{\theta}) \mathbf{x}$$

The following one-to-one relation exists between  $\boldsymbol{\theta}$  and  $\boldsymbol{\eta}$ :

$$\boldsymbol{\eta} = \partial_{\boldsymbol{\theta}} \psi(\boldsymbol{\theta}) \quad (4)$$

[Theorem 2] The marginalized probability distribution  $\Pi \circ q(\mathbf{x})$  of  $q(\mathbf{x})$  is an  $m$ -projection of  $q(\mathbf{x})$  on  $M_D$ .

(Proof) Let us consider an  $m$ -projection of  $q(\mathbf{x})$  to  $M_D$ . From Theorem 1, we can differentiate  $D[q(\mathbf{x}); p(\mathbf{x}; \boldsymbol{\theta})]$  with respect to  $\boldsymbol{\theta}$ . Using the results of Eq. (4), the following is obtained:

$$\partial_{\boldsymbol{\theta}} D[q(\mathbf{x}); p(\mathbf{x}; \boldsymbol{\theta})] = \boldsymbol{\eta} - \sum_{\mathbf{x}} q(\mathbf{x}) \mathbf{x}$$

Consequently,  $\boldsymbol{\eta}^* = \sum_{\mathbf{x}} q(\mathbf{x}) \mathbf{x}$ , where  $\boldsymbol{\eta}^*$  is the projected  $\boldsymbol{\eta}$  coordinate, gives the  $m$ -projection. This indicates that an  $m$ -projection is expressed by the independent expected value of each component of  $\mathbf{x}$  based on  $q(\mathbf{x})$ , which proves the theorem.

Furthermore, MPM decoding can be written as follows.

$$\hat{\mathbf{x}} = \text{sgn}(\boldsymbol{\eta}^*)$$

where  $\text{sgn}(\cdot)$  is applied independently to each bit.

#### 2.4. Information geometrical view of turbo decoding

Turbo decoding does not use two observed parity check words simultaneously as  $p(\mathbf{x}|\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$ , but uses  $p(\tilde{\mathbf{x}}|\tilde{\mathbf{y}}_1|\mathbf{x})$ ,  $p(\mathbf{x}|\tilde{\mathbf{y}}_2|\mathbf{x})$  in which  $\tilde{\mathbf{x}}$  and only one of the observed parity check words is taken into account. The following gives the form of  $p(\tilde{\mathbf{x}}|\tilde{\mathbf{y}}_r|\mathbf{x})$ ,  $r = 1, 2$ :

$$p(\tilde{\mathbf{x}}|\tilde{\mathbf{y}}_r|\mathbf{x}) = \exp(c_0(\mathbf{x}) + c_r(\mathbf{x}) - (N + L)\psi(\beta))$$

The following is obtained for these distributions when  $\omega(\mathbf{x}; \boldsymbol{\xi}) \in M_D$  is taken to be the prior distribution of  $\mathbf{x}$ , and a posterior distribution is given as follows:

$$\begin{aligned} p_r(\mathbf{x}; \boldsymbol{\xi}) &= p_r(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_r; \boldsymbol{\xi}) \\ &= \frac{p_r(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_r|\mathbf{x})\omega(\mathbf{x}; \boldsymbol{\xi})}{\sum_{\mathbf{x}} p_r(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_r|\mathbf{x})\omega(\mathbf{x}; \boldsymbol{\xi})} \\ &= \exp(c_0(\mathbf{x}) + c_r(\mathbf{x}) + \boldsymbol{\xi} \cdot \mathbf{x} - \varphi_r(\boldsymbol{\xi})) \end{aligned}$$

$\varphi_r(\boldsymbol{\xi})$  is the normalization factor. Here, we assume that the  $m$ -projection from  $p_r(\mathbf{x}; \boldsymbol{\xi})$ ,  $r = 1, 2$ , to  $M_D$  is tractable in polynomial time. The turbo decoding varies iteratively, and the  $m$ -projection of  $p(\mathbf{x}|\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$  is approximated.

In this section, the information geometrical view of turbo decoding is given. First, we define three important manifolds:

$$\begin{aligned} M_0 &= \left\{ p_0(\mathbf{x}; \boldsymbol{\xi}) = \exp(c_0(\mathbf{x}) + \boldsymbol{\xi} \cdot \mathbf{x} - \varphi_0(\boldsymbol{\xi})) \mid \boldsymbol{\xi} \in \mathcal{R}^N \right\} \end{aligned}$$

$$M_1 = \left\{ p_1(\mathbf{x}; \boldsymbol{\xi}) \mid \boldsymbol{\xi} \in \mathcal{R}^N \right\}$$

$$M_2 = \left\{ p_2(\mathbf{x}; \boldsymbol{\xi}) \mid \boldsymbol{\xi} \in \mathcal{R}^N \right\}$$

Here,  $\boldsymbol{\xi}$  defines the coordinate system of each manifold. Since  $c_0(\mathbf{x}) = \beta \tilde{\mathbf{x}} \cdot \mathbf{x}$ , the condition  $p(\mathbf{x}; \boldsymbol{\theta}') = p_0(\mathbf{x}; \boldsymbol{\theta})$  holds for  $p(\mathbf{x}; \boldsymbol{\theta}') \in M_D$ , if  $\boldsymbol{\theta}' = \boldsymbol{\theta} + \beta \tilde{\mathbf{x}}$  is assumed. Consequently,  $M_0$  is equivalent to  $M_D$ . Let  $\pi_{M_0} \circ q(\mathbf{x})$  be a coordinate  $\boldsymbol{\xi}$  defined by the  $m$ -projection of  $q(\mathbf{x})$  on  $M_0$ :

$$\pi_{M_0} \circ q(\mathbf{x}) = \underset{\boldsymbol{\xi} \in \mathcal{R}^N}{\text{argmin}} D[q(\mathbf{x}); p_0(\mathbf{x}; \boldsymbol{\xi})]$$

The condition  $\tau_{M_0} \circ q(\mathbf{x}) = \pi_{M_0} \circ q(\mathbf{x}) + \beta \tilde{\mathbf{x}}$  holds. The turbo decoding is written in the following manner with  $\pi_{M_0}$  (Fig. 3).

[Expression of Turbo Decoding as Information Geometry]

1. Set  $\xi_1^t = 0$  for  $t = 0$ , and  $t = 1$ .
2. Calculate  $\pi_{M_0} \circ p_2(\mathbf{x}; \xi_1^t)$ , which is the projection of  $p_2(\mathbf{x}; \xi_1^t)$  to  $M_0$ , and update  $\xi_2^{t+1}$  as follows:

$$\xi_2^{t+1} = \pi_{M_0} \circ p_2(\mathbf{x}; \xi_1^t) - \xi_1^t \quad (5)$$

3. Calculate  $\pi_{M_0} \circ p_1(\mathbf{x}; \xi_2^{t+1})$ , which is the projection of  $p_1(\mathbf{x}; \xi_2^{t+1})$  to  $M_0$ , and update  $\xi_1^{t+1}$  as follows:

$$\xi_1^{t+1} = \pi_{M_0} \circ p_1(\mathbf{x}; \xi_2^{t+1}) - \xi_2^{t+1} \quad (6)$$

4. Return to step 2 if  $\pi_{M_0} \circ p_1(\mathbf{x}; \xi_2^{t+1}) \neq \tau_{M_0} \circ p_2(\mathbf{x}; \xi_1^{t+1})$ .

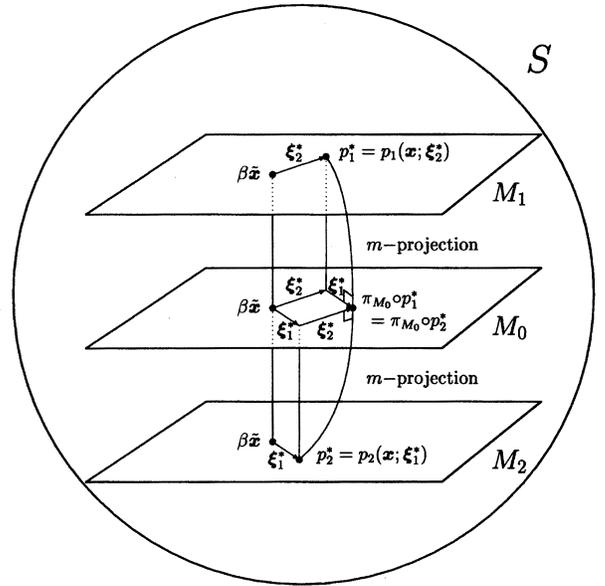


Fig. 3. Information geometry of turbo decoding.

$L_1 \mathbf{x}^{(l)}, L_2 \mathbf{x}^{(l)}$  in Eqs. (1) and (2), that is, the expressions  $F((\mathbf{l}\mathbf{x} + \xi_1), \mathbf{l}y_1), F((\mathbf{l}\mathbf{x} + \xi_2), \mathbf{l}y_2)$  correspond to  $\tau_{M_0} \circ p_2(\mathbf{x}; \xi_1^*), \pi_{M_0} \circ p_1(\mathbf{x}; \xi_2^{*+1})$  in Eqs. (5) and (6). In addition,  $\xi_1, \xi_2$  in Eqs. (1) and (2) correspond to  $\xi_1^*, \xi_2^{*+1}$  in Eqs. (5) and (6).

### 3. Properties of Turbo Decoding

#### 3.1. Properties of stationary points

Let us assume that  $\xi_1^*, \xi_2^*$  are the convergence points of turbo decoding. The final result is the  $M_0$  coordinate at which  $\tau_{M_0} \circ p_1(\mathbf{x}; \xi_2^*) = \pi_{M_0} \circ p_2(\mathbf{x}; \xi_1^*)$ . Let us define the point as  $\theta^*$ . First, the following is obtained based on convergence conditions.

$$\Pi \circ p_1(\mathbf{x}; \xi_2^*) = \Pi \circ p_2(\mathbf{x}; \xi_1^*) = p_0(\mathbf{x}; \theta^*)$$

Also, the following is obtained based on  $\tau_{M_0} \circ p_1(\mathbf{x}; \xi_2^*) = \pi_{M_0} \circ p_2(\mathbf{x}; \xi_1^*)$  and on steps 2 and 3 of the algorithm:

$$\theta^* = \xi_1^* + \xi_2^*$$

In turbo decoding, the results of true MPM decoding are approximated as  $\theta^* = \xi_1^* + \xi_2^*$ :

$$p_0(\mathbf{x}; \theta^*) = \exp(c_0(\mathbf{x}) + \xi_1^* \cdot \mathbf{x} + \xi_2^* \cdot \mathbf{x} - \varphi_0(\theta^*)) \quad (7)$$

Intuitively,  $\xi_2$  in Eq. (7) is substituted by  $c_2(\mathbf{x})$  in step 2,  $\xi_1$  in Eq. (7) is substituted by  $c_1(\mathbf{x})$  in step 3 of turbo decoding, and  $\xi_1^*$  is determined. Consequently, the influences of  $c_1(\mathbf{x})$  and  $c_2(\mathbf{x})$  are expressed by  $\xi_1^*$  and  $\xi_2^*$ , respectively, but their influences usually cannot be linearly separated on  $M_0$ .

Let us define the  $\xi_1$  and  $\xi_2$  that satisfy the equation below as  $\xi_1(\theta)$  and  $\xi_2(\theta)$ , respectively:

$$\Pi \circ p_1(\mathbf{x}; \xi_2) = \Pi \circ p_2(\mathbf{x}; \xi_1) = p_0(\mathbf{x}; \theta)$$

Also, let us define the  $m$ - and  $e$ -flat manifolds that connect  $p_0(\mathbf{x}; \theta), p_1(\mathbf{x}; \xi_2(\theta)), p_2(\mathbf{x}; \xi_1(\theta))$  as  $M(\theta)$  and  $E(\theta)$ , respectively:

$$M(\theta) = \left\{ p(\mathbf{x}) \mid \sum_{\mathbf{x}} p(\mathbf{x}) \mathbf{x} = \sum_{\mathbf{x}} p_0(\mathbf{x}; \theta) \mathbf{x} \right\}$$

$$E(\theta) = \left\{ p = C p_0^{t_0} p_1^{t_1} p_2^{t_2} \mid \sum_{r=0}^2 t_r = 1 \right\}$$

For every  $p(\mathbf{x}) \in M(\theta)$ , its  $m$ -projection to  $M_0$  is  $p_0(\mathbf{x}; \theta)$ .

[Theorem 3] At a stationary point, the three distributions  $p_0^*, p_1^*, p_2^*$  are included in  $M(\theta^*)$ , and the four distributions  $p_0^*, p_1^*, p_2^*, p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$  are included in  $E(\theta^*)$  (Fig. 4).

(Proof) Based on the definition, the fact that  $p_0^*, p_1^*, p_2^* \in (M(\theta^*), E(\theta^*))$  is proven. The fact that the four distributions  $p_0^*, p_1^*, p_2^*, p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) \in E(\theta^*)$  is proven from the following result:

$$\begin{aligned} & C \frac{p_1(\mathbf{x}; \xi_2^*) p_2(\mathbf{x}; \xi_1^*)}{p_0(\mathbf{x}; \theta^*)} \\ &= C \exp(c_0(\mathbf{x}) + c_1(\mathbf{x}) + c_2(\mathbf{x})) \\ &= p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) \end{aligned}$$

by setting  $t_0 = -1, t_1 = t_2 = 1$  and using  $\theta^* = \xi_1^* + \xi_2^*$ .

$M(\theta_{MPM}^*)$  includes  $p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$  if a true solution for MPM decoding is  $\theta_{MPM}^*$ . Generally, however, the solution  $\theta^*$  for turbo decoding and  $\theta_{MPM}^*$  do not agree with each other. Consequently,  $p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2) \notin M(\theta^*)$ , but is included in  $E(\theta^*)$ . Since  $e$ -flatness and  $m$ -flatness generally do not agree with each other, there is a discrepancy between the manifolds  $E(\theta^*)$  and  $M(\theta^*)$ . Turbo decoding obtains the approximation by substituting  $M(\theta^*)$  for  $E(\theta^*)$ . Similar structures exist in other statistical physics techniques [9–12].

#### 3.2. Stability of stationary point

Let us assume that the convergence points of turbo decoding are  $\xi_1^*, \xi_2^*, \theta^* = \xi_1^* + \xi_2^*$ .

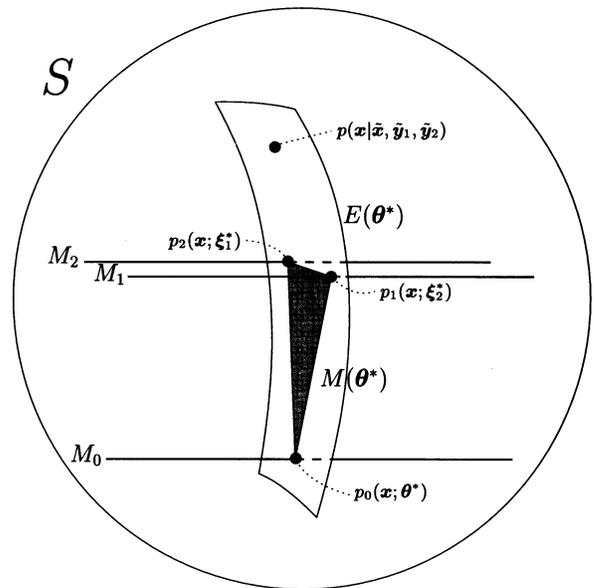


Fig. 4. Information geometrical view of turbo codes.

Now we define  $G_0(\boldsymbol{\theta})$ ,  $G_1(\boldsymbol{\xi})$ ,  $G_2(\boldsymbol{\xi})$  as the Fisher information matrices of  $p_0(\mathbf{x}; \boldsymbol{\theta})$ ,  $p_1(\mathbf{x}; \boldsymbol{\xi})$ ,  $p_2(\mathbf{x}; \boldsymbol{\xi})$ , respectively;  $I_N$  is a unit matrix, and  $\boldsymbol{\eta}_0(\boldsymbol{\theta})$ ,  $\boldsymbol{\eta}_1(\boldsymbol{\xi})$ ,  $\boldsymbol{\eta}_2(\boldsymbol{\xi})$  are the expectation parameters of the corresponding distributions. The following condition holds for the convergence points  $\boldsymbol{\eta}_0(\boldsymbol{\theta}^*)$ ,  $\boldsymbol{\eta}_1(\boldsymbol{\xi}_2^*) = \boldsymbol{\eta}_2(\boldsymbol{\xi}_1^*)$ . The following also holds for each of these ( $r = 0, 1, 2$ ):

$$G_r(\boldsymbol{\theta}) = \partial_{\boldsymbol{\theta}\boldsymbol{\theta}'} \varphi_r(\boldsymbol{\theta}) = \partial_{\boldsymbol{\theta}} \boldsymbol{\eta}_r(\boldsymbol{\theta})$$

To see the stability, we add a sufficiently small vector  $\boldsymbol{\delta}$  to  $\boldsymbol{\xi}_1^*$  as  $\boldsymbol{\xi}_1 = \boldsymbol{\xi}_1^* + \boldsymbol{\delta}$ , and set it as the initial value of the algorithm, then the turbo decoding algorithm is applied once. Setting  $\boldsymbol{\xi}_1' = \boldsymbol{\xi}_1^* + \boldsymbol{\delta}'$  to be the parameter after a single cycle of turbo decoding, a linear stability analysis is carried out. From step 2, the following is obtained:

$$\boldsymbol{\eta}_0(\boldsymbol{\theta}^* + \Delta\boldsymbol{\theta}) = \boldsymbol{\eta}_2(\boldsymbol{\xi}_1^* + \boldsymbol{\delta})$$

When this expression is expanded, the result is

$$\begin{aligned} \boldsymbol{\eta}_0(\boldsymbol{\theta}^*) + G_0(\boldsymbol{\theta}^*)\Delta\boldsymbol{\theta} &= \boldsymbol{\eta}_1(\boldsymbol{\xi}_1^*) + G_2(\boldsymbol{\xi}_1^*)\boldsymbol{\delta} \\ \Delta\boldsymbol{\theta} &= G_0(\boldsymbol{\theta}^*)^{-1}G_2(\boldsymbol{\xi}_1^*)\boldsymbol{\delta} \end{aligned}$$

Also,  $\boldsymbol{\xi}_2$  in step 2 becomes

$$\boldsymbol{\xi}_2 = \boldsymbol{\xi}_2^* + (G_0(\boldsymbol{\theta}^*)^{-1}G_2(\boldsymbol{\xi}_1^*) - I_N) \boldsymbol{\delta}$$

Step 3 is treated in the same manner. As a result,  $\boldsymbol{\delta}$  is updated as

$$\boldsymbol{\delta}' = \mathcal{T} \boldsymbol{\delta}$$

$$\mathcal{T} = (G_0(\boldsymbol{\theta}^*)^{-1}G_1(\boldsymbol{\xi}_2^*) - I_N)(G_0(\boldsymbol{\theta}^*)^{-1}G_2(\boldsymbol{\xi}_1^*) - I_N)$$

which gives the linearized approximation of turbo decoding.

[Theorem 4] Let  $\lambda_i$  be the eigenvalue of  $\mathcal{T}$ . The stationary point will be stable if  $|\lambda_i| < 1$  holds for every  $i$ .

This result agrees with that of Ref. 2.

### 3.3. Cost function and stationary points

Let us consider the following function, assuming  $\boldsymbol{\theta} = \boldsymbol{\xi}_1 + \boldsymbol{\xi}_2$ :

$$\mathcal{F}(\boldsymbol{\xi}_1, \boldsymbol{\xi}_2) = \varphi_0(\boldsymbol{\theta}) - \varphi_1(\boldsymbol{\xi}_2) - \varphi_2(\boldsymbol{\xi}_1)$$

[Theorem 5] The stationary points  $\boldsymbol{\xi}_1^*$ ,  $\boldsymbol{\xi}_2^*$  of turbo decoding are the critical points of  $\mathcal{F}$ .

(Proof) With direct differentiation, the following is obtained:

$$\begin{aligned} \partial_{\boldsymbol{\xi}_1} \mathcal{F} &= \partial_{\boldsymbol{\theta}} \varphi_0(\boldsymbol{\theta}) - \partial_{\boldsymbol{\xi}_1} \varphi_2(\boldsymbol{\xi}_1) = \boldsymbol{\eta}_0(\boldsymbol{\theta}) - \boldsymbol{\eta}_2(\boldsymbol{\xi}_1) \\ \partial_{\boldsymbol{\xi}_2} \mathcal{F} &= \partial_{\boldsymbol{\theta}} \varphi_0(\boldsymbol{\theta}) - \partial_{\boldsymbol{\xi}_2} \varphi_1(\boldsymbol{\xi}_2) = \boldsymbol{\eta}_0(\boldsymbol{\theta}) - \boldsymbol{\eta}_1(\boldsymbol{\xi}_2) \end{aligned}$$

The above differential is 0 because  $\boldsymbol{\eta}_0(\boldsymbol{\theta}^*) = \boldsymbol{\eta}_2(\boldsymbol{\xi}_1^*) = \boldsymbol{\eta}_1(\boldsymbol{\xi}_2^*)$  at the equilibrium.

The turbo decoding algorithm can be approximated with

$$\begin{pmatrix} \boldsymbol{\xi}_1^{t+1} \\ \boldsymbol{\xi}_2^{t+1} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\xi}_1^t \\ \boldsymbol{\xi}_2^t \end{pmatrix} \simeq \begin{pmatrix} O & G_0(\boldsymbol{\theta})^{-1} \\ G_0(\boldsymbol{\theta})^{-1} & O \end{pmatrix} \begin{pmatrix} \partial_{\boldsymbol{\xi}_1} \mathcal{F} \\ \partial_{\boldsymbol{\xi}_2} \mathcal{F} \end{pmatrix}$$

when the parameter changes by a very small value. The Hessian of  $\mathcal{F}$  is calculated as

$$\mathcal{H} = \begin{pmatrix} \partial_{\boldsymbol{\xi}_1 \boldsymbol{\xi}_1} \mathcal{F} & \partial_{\boldsymbol{\xi}_1 \boldsymbol{\xi}_2} \mathcal{F} \\ \partial_{\boldsymbol{\xi}_2 \boldsymbol{\xi}_1} \mathcal{F} & \partial_{\boldsymbol{\xi}_2 \boldsymbol{\xi}_2} \mathcal{F} \end{pmatrix} = \begin{pmatrix} G_0 - G_1 & G_0 \\ G_0 & G_0 - G_2 \end{pmatrix}$$

Here, the variables are changed to  $\boldsymbol{\theta} = \boldsymbol{\xi}_1 + \boldsymbol{\xi}_2$ ,  $\boldsymbol{\nu} = \boldsymbol{\xi}_1 - \boldsymbol{\xi}_2$ , yielding

$$\begin{aligned} &\begin{pmatrix} \partial_{\boldsymbol{\theta}\boldsymbol{\theta}} \mathcal{F} & \partial_{\boldsymbol{\theta}\boldsymbol{\nu}} \mathcal{F} \\ \partial_{\boldsymbol{\nu}\boldsymbol{\theta}} \mathcal{F} & \partial_{\boldsymbol{\nu}\boldsymbol{\nu}} \mathcal{F} \end{pmatrix} \\ &= \frac{1}{4} \begin{pmatrix} 4G_0(\boldsymbol{\theta}) - (G_1 + G_2) & (G_1 - G_2) \\ (G_1 - G_2) & -(G_1 + G_2) \end{pmatrix} \end{aligned}$$

from which we can see that  $\partial_{\boldsymbol{\theta}\boldsymbol{\theta}} \mathcal{F}$  is probably positive definite and that  $\partial_{\boldsymbol{\nu}\boldsymbol{\nu}} \mathcal{F}$  is always negative. Consequently, the stationary points of turbo decoding are often believed to be a saddle point.

### 3.4. Difference between turbo decoding results and MPM solution

Theorem 3 shows that the difference between MPM decoding and turbo decoding is the difference between  $M(\boldsymbol{\theta})$  and  $E(\boldsymbol{\theta})$ . Based on this result, we evaluate the difference between the true MPM solution and the turbo decoding solution by the perturbation analysis. First,  $\boldsymbol{\theta} = (\boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^N)^T$ ,  $\mathbf{v} = (v^1, v^2)^T$ , and  $\mathbf{c}(\mathbf{x}) \stackrel{\text{def}}{=} (c_1(\mathbf{x}), c_2(\mathbf{x}))^T$  are used to define  $p(\mathbf{x}; \boldsymbol{\theta}, \mathbf{v})$  in the following manner:

$$\begin{aligned} p(\mathbf{x}; \boldsymbol{\theta}, \mathbf{v}) &= \exp(c_0(\mathbf{x}) + \boldsymbol{\theta} \cdot \mathbf{x} + \mathbf{v} \cdot \mathbf{c}(\mathbf{x}) - \varphi(\boldsymbol{\theta}, \mathbf{v})) \\ \varphi(\boldsymbol{\theta}, \mathbf{v}) &= \ln \sum_{\mathbf{x}} \exp(c_0(\mathbf{x}) + \boldsymbol{\theta} \cdot \mathbf{x} + \mathbf{v} \cdot \mathbf{c}(\mathbf{x})) \end{aligned}$$

This distribution includes  $p_0(\mathbf{x}; \boldsymbol{\theta})(\mathbf{v} = 0)$ ,  $p(\mathbf{x}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$  ( $\boldsymbol{\theta} = \mathbf{0}$ ,  $\mathbf{v} = (1, 1)^T$ ),  $p_r(\mathbf{x}; \boldsymbol{\xi})$  ( $\boldsymbol{\theta} = \boldsymbol{\xi}$ ,  $\mathbf{v} = \mathbf{e}_r$ ) as special cases. Here,  $\mathbf{e}_r$  is

$$\mathbf{e}_1 = (1, 0)^T, \quad \mathbf{e}_2 = (0, 1)^T$$

The expectation parameter  $\boldsymbol{\eta}(\boldsymbol{\theta}, \mathbf{v})$  is defined as follows:

$$\boldsymbol{\eta}(\boldsymbol{\theta}, \mathbf{v}) = \partial_{\boldsymbol{\theta}} \varphi(\boldsymbol{\theta}, \mathbf{v}) = \sum_{\mathbf{x}} p(\mathbf{x}; \boldsymbol{\theta}, \mathbf{v}) \mathbf{x}$$

This defines a submanifold  $M(\boldsymbol{\theta}^*)$  in which the expectation parameters of all distributions are equal to  $\boldsymbol{\eta}(\boldsymbol{\theta}^*)$ . In other words, we have a set of distributions for which the condition  $\boldsymbol{\eta}(\boldsymbol{\theta}, \mathbf{v}) = \boldsymbol{\eta}(\boldsymbol{\theta}^*)$  is satisfied. Hereafter, we use the following notations  $\boldsymbol{\eta}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \boldsymbol{\eta}(\boldsymbol{\theta}, \mathbf{0})$  and  $\boldsymbol{\eta}^* \stackrel{\text{def}}{=} \boldsymbol{\eta}(\boldsymbol{\theta}^*)$ .

The dependency of  $\boldsymbol{\theta}$  on  $\mathbf{v}$  will now be considered using the perturbation analysis under the restriction that  $p(\mathbf{x}; \boldsymbol{\theta}, \mathbf{v}) \in M(\boldsymbol{\theta}^*)$ . In the perturbation analysis, the contribution of  $\{v^r c_r(\mathbf{x})\}$  to  $p(\mathbf{x}; \boldsymbol{\theta}, \mathbf{v})$  is assumed to be small, we take the Taylor series of  $p(\mathbf{x}; \boldsymbol{\theta}, \mathbf{v})$  with respect to  $\{v^r\}$  up to  $O(\|\mathbf{v}\|^2)$ , and the effect of  $\{v^r c_r(\mathbf{x})\}$  is evaluated. The Taylor expansion gives

$$\begin{aligned} \eta_i(\boldsymbol{\theta}, \mathbf{v}) &= \eta_i^* + \sum_j \partial_j \eta_i^* \Delta \theta^j + \sum_r \partial_r \eta_i^* v^r \\ &+ \frac{1}{2} \sum_{r,s} \partial_r \partial_s \eta_i^* v^r v^s + \sum_{j,r} \partial_r \partial_j \eta_i^* v^r \Delta \theta^j \\ &+ \frac{1}{2} \sum_{k,l} \partial_k \partial_l \eta_i^* \Delta \theta^k \Delta \theta^l + O(\|\mathbf{v}\|^3) + O(\|\Delta \boldsymbol{\theta}\|^3) \end{aligned}$$

where  $\{i, j, k, l\}$  are the subscripts of  $\boldsymbol{\theta}$ ,  $\{r, s\}$  are the subscripts of  $\mathbf{v}$ , and  $\Delta \boldsymbol{\theta} \stackrel{\text{def}}{=} \boldsymbol{\theta} - \boldsymbol{\theta}^*$ . From the condition that the distribution must be on  $M(\boldsymbol{\theta}^*)$ ,  $\eta_i(\boldsymbol{\theta}, \mathbf{v}) = \eta_i^*$  holds. Also,  $\{g_{ij}\}$ , which is the Fisher information matrix of  $p(\mathbf{x}; \boldsymbol{\theta}^*, \mathbf{0})$ , is a diagonal matrix. Using these results, the following equation is derived:

$$\begin{aligned} \Delta \theta^i &= -g^{ii} \left[ \sum_r \partial_r \eta_i^* v^r - \frac{1}{2} \sum_{r,s} \partial_r \partial_s \eta_i^* v^r v^s \right. \\ &\quad \left. - \frac{1}{2} \sum_{k,l} \partial_k \partial_l \eta_i^* \Delta \theta^k \Delta \theta^l - \sum_{k,r} \partial_r \partial_k \eta_i^* v^r \Delta \theta^k \right] \\ &+ O(\|\mathbf{v}\|^3) + O(\|\Delta \boldsymbol{\theta}\|^3) \end{aligned}$$

where  $g^{ii} = 1/g_{ii}$ . By ignoring third- and higher-order terms of  $v^r \Delta \theta^l$  is rewritten as follows:

$$\begin{aligned} \Delta \theta^i &\simeq -g^{ii} \sum_r A_{ir} v^r - \frac{g^{ii}}{2} \times \\ &\sum_{r,s} \left( \partial_r - \sum_k g^{kk} A_{kr} \partial_k \right) \left( \partial_s - \sum_j g^{jj} A_{js} \partial_j \right) \eta_i^* v^r v^s \end{aligned} \quad (8)$$

where  $A_{ir} = \partial_r \eta_i^*$ .

When  $\mathbf{v} = \mathbf{e}_1$ ,  $p(\mathbf{x}; \boldsymbol{\theta}, \mathbf{v})$  is restricted to  $M(\boldsymbol{\theta}^*)$ , so  $\boldsymbol{\theta} = \boldsymbol{\xi}_2^*$  and  $\Delta \boldsymbol{\theta} = \boldsymbol{\xi}_2^* - \boldsymbol{\theta}^* = -\boldsymbol{\xi}_1^*$ . Based on the same argument,  $\Delta \boldsymbol{\theta} = -\boldsymbol{\xi}_2^*$  when  $\mathbf{v} = \mathbf{e}_2$ . The following expression is derived from Eq. (8):

$$-\xi_r^{i,*} \simeq -g^{ii} A_{ir}$$

$$-\frac{g^{ii}}{2} \left( \partial_r - \sum_k g^{kk} A_{kr} \partial_k \right) \left( \partial_r - \sum_j g^{jj} A_{jr} \partial_j \right) \eta_i^* \quad (9)$$

On the other hand, if we assume that  $\mathbf{v} = \sum_r \mathbf{e}_r$  and define  $\bar{\boldsymbol{\theta}}$  as the parameter that satisfies the condition, generally,  $\bar{\boldsymbol{\theta}} \neq \boldsymbol{\theta}^*$ . This indicates that  $p(\mathbf{x}; \bar{\boldsymbol{\theta}}, \tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2)$  is not necessarily included in  $M(\boldsymbol{\theta}^*)$ . It follows from Eq. (8) that

$$\begin{aligned} \bar{\theta}^i - \theta^{i,*} &\simeq -g^{ii} \sum_r A_{ir} \\ &- \frac{g^{ii}}{2} \sum_r \left( \partial_r - \sum_k g^{kk} A_{kr} \partial_k \right) \left( \partial_r - \sum_j g^{jj} A_{jr} \partial_j \right) \eta_i^* \end{aligned}$$

From  $\theta^{i,*} = \xi_1^{i,*} + \xi_2^{i,*}$ , and based on the results of Eq. (9), we can conclude that

$$\begin{aligned} \bar{\theta}^i &\simeq \\ &- \frac{g^{ii}}{2} \sum_{r \neq s} \left( \partial_r - \sum_k g^{kk} A_{kr} \partial_k \right) \left( \partial_s - \sum_j g^{jj} A_{js} \partial_j \right) \eta_i^* \end{aligned}$$

The difference between a turbo solution and an MPM solution is evaluated if the difference between the parameters on  $M_0$  is evaluated. If the MPM solution can be written as  $\boldsymbol{\eta}_{MPM}^*$ , the result is

$$\begin{aligned} \eta_{i,MPM}^* &\simeq \eta_i^* \\ &+ \frac{1}{2} \sum_{r \neq s} \left( \partial_r - \sum_k g^{kk} A_{kr} \partial_k \right) \left( \partial_s - \sum_j g^{jj} A_{js} \partial_j \right) \eta_i^* \end{aligned}$$

The following theorem is derived from this result.

[Theorem 6] Assuming that  $\boldsymbol{\eta}^*$  is the expectation of  $x_i$  based on a turbo decoding solution, and  $\boldsymbol{\eta}_{MPM}^*$  is the expectation based on an MPM decoding solution, the difference between these is approximated as follows:

$$\begin{aligned} \boldsymbol{\eta}_{MPM}^* - \boldsymbol{\eta}^* &\simeq \\ &\frac{1}{2} \sum_{r \neq s} \left( \partial_r - \sum_k g^{kk} A_{kr} \partial_k \right) \left( \partial_s - \sum_j g^{jj} A_{js} \partial_j \right) \boldsymbol{\eta}^* \end{aligned}$$

The decoding error given by the above formula is related to the embedded  $e$ -curvature of the manifold  $M(\boldsymbol{\theta})$ .

## 4. Conclusions

In this paper, we have elucidated the mathematical structure of turbo codes from an information geometrical viewpoint. We believe that our results provide a mathematical framework for analyzing turbo codes and that more properties will be elucidated based on this framework.

This problem of turbo decoding can be regarded in a more general sense, as a problem of approximating the marginalization of a probability distribution which has a form shown in Eq. (3). The solution is approximated by dividing the whole problem into partial problems and partial information is integrated through iterative algorithm.

The same structure is found in low-density parity-check codes, the Bethe approximation of statistical physics, and the belief propagation of loopy Bayesian networks. Some differences exist in relation to the details of the algorithm, and completely identical results have not been obtained concerning solution stability and other properties, but since the basic structures are very similar, we believe it is not difficult to build the information geometrical framework for related problems. We have already extended the framework to the case of low-density parity-check codes [13]. Extensions for the Bethe approximation and belief propagation techniques are our future works.

**Acknowledgments.** We gratefully acknowledge Yoshiyuki Kabashima and Motohiko Isaka for their advice regarding this study. This work was supported in part by MEXT, Japan under a Grant-in-Aid for Scientific Research (16700227, 14084208).

## REFERENCES

1. Berrou C, Glavieux A, Thitimajshima P. Near Shannon limit error-correcting coding and decoding: Turbo-codes. Proc IEEE Int Conf on Communications, p 1064–1070, Geneva, 1993.
2. Richardson T. The geometry of turbo-decoding dynamics. IEEE Trans Inf Theory 2000;46:9–23.
3. Gallager RG. Low density parity check codes. IRE Trans Inf Theory 1962;8:21–28.
4. MacKay DJC. Good error-correcting codes based on very sparse matrices. IEEE Trans Inf Theory 1999;45:399–431.
5. McEliece RJ, MacKay DJC, Cheng JF. Turbo decoding as an instance of Pearl’s “belief propagation” algorithm. IEEE J Sel Areas Commun 1998;16:140–152.
6. Kabashima Y, Saad D. The TAP approach to intensive and extensive connectivity systems. In: Oppor M, Saad D (editors). Advanced mean field methods—Theory and practice. MIT Press; 2001. p 65–84.
7. Amari S. Differential-geometrical methods in statistics. Lecture Notes in Statistics Vol. 28. Springer-Verlag; 1985.
8. Amari S, Nagaoka H. Methods of information geometry. AMS and Oxford University Press; 2000.
9. Kappen HJ, Wiegnerck WJ. Mean field theory graphical models. In: Oppor M, Saad D (editors). Advanced mean field methods—Theory and practice. MIT Press; 2001. p 37–49.
10. Amari S, Ikeda S, Shimokawa H. Information geometry and mean field approximation: The  $\alpha$ -projection approach. In: Oppor M, Saad D (editors). Advanced mean field methods—Theory and practice. MIT Press; 2001. p 241–257.
11. Tanaka T. Information geometry of mean-field approximation. Neural Computation 2000;12:1951–1968.
12. Tanaka T. Information geometry of mean-field approximation. In: Oppor M, Saad D (editors). Advanced mean field methods—Theory and practice. MIT Press; 2001. p 259–273.
13. Ikeda S, Tanaka T, Amari S. Information geometry of turbo codes and low-density parity-check codes. IEEE Trans Inf Theory 2004;50:1097–1114.

## AUTHORS (from left to right)



**Shiro Ikeda** (member) received his B.E., M.E., and D.Eng. degrees in information physics from the University of Tokyo in 1991, 1993, and 1996. From 1996 to 2001, he was with RIKEN BSI, first as a special postdoctoral researcher of RIKEN and later as a researcher of JST. He was an associate professor at Kyushu Institute of Technology from 2001 to 2003, and has been an associate professor at the Institute of Statistical Mathematics since 2003. His research interests are in the areas of statistical signal processing, learning theory, and information geometry. He received the Best Research Award and Best Paper Award from the Japan Neural Network Society in 1999 and 2001, respectively.

**Toshiyuki Tanaka** (member) received his B.E., M.E., and D.Eng. degrees in electronics engineering from the University of Tokyo in 1988, 1990, and 1993. In 1993, he joined the Department of Electronics and Information Engineering, Tokyo Metropolitan University, and is currently an associate professor there. His research interests are in the interdisciplinary areas of information and communication theory, learning theory, information geometry, and statistical mechanics. He received the DoCoMo Mobile Science Award in 2002.

**Shun-ichi Amari** (fellow) graduated from the University of Tokyo in 1958, majoring in mathematical engineering, and received a D.Eng. degree from the University of Tokyo in 1963. He was an associate professor at Kyushu University, an associate professor and professor in the Department of Mathematical Engineering and Information Physics, University of Tokyo, and is now professor emeritus at the University of Tokyo. He is the Director of RIKEN BSI. He has been engaged in research in wide areas of mathematical engineering and applied mathematics, such as topological network theory, differential geometry of continuum mechanics, pattern recognition, mathematical foundations of neural networks, and information geometry. He served as President of the International Neural Network Society, Council member of the Bernoulli Society for Mathematical Statistics and Probability Theory, and is President-Elect of IEICE. He was founding Co-Editor-in-Chief of *Neural Networks*. He has been awarded the Japan Academy Award, IEEE Neural Networks Pioneer Award, IEEE Emanuel R. Poire Award, Neurocomputing best paper award, IEEE Signal Processing Society best paper award, and NEC C&C Prize, among many others.