

情報幾何学に基づく確率伝搬法の解析

池田 思朗¹・田中 利幸²・甘利 俊一³

(受付 2004 年 2 月 27 日 ; 改訂 2004 年 6 月 29 日)

要 旨

1980年代後半 Pearl が提案した確率伝搬法は、大規模なグラフィカルモデルに対する確率推論のための計算手法である。同等の手法は統計物理学、統計学、誤り訂正符号の復号法などにも存在し、広く用いられている。確率伝搬法は木の構造のグラフに対してはグラフの大きさに比例した計算量で厳密解が得られる。しかしループを持つグラフに対しては繰り返し計算の収束性、および得られた結果の近似精度ともに理論的には十分理解されていなかった。一方で確率伝搬法は実用上有効な手法であり、その性質を理論的に明らかにすることは重要である。本研究では情報幾何学に基づく枠組みにより確率伝搬法を表現し、収束性や近似精度を議論する。

キーワード： 確率伝搬法，情報幾何学，グラフィカルモデル。

1. はじめに

グラフィカルモデル(Lauritzen and Spiegelhalter (1988), Jordan (1999))では、複数の確率変数の同時分布をグラフによって表現する。グラフによって表現された複数の確率変数の一部のみが観測されたとき、その条件付確率分布から観測されていない確率変数の値を推論する問題を考える。この問題は確率推論と呼ばれ、人工知能、統計物理学、情報理論など様々な分野で重要である。

Pearl (1988)はこの問題に対し、確率伝搬法(Belief Propagation)と呼ばれる簡便な繰り返し計算法を提案した。確率伝搬法は条件付き分布から各確率変数の周辺分布を求め、確率推論を得る。木のグラフに対しては収束が保証され、グラフの辺の数に比例した計算量で正しい推論結果が得られることが分かっているが、グラフがループを持つ場合には必ずしも収束しないこと、収束した場合でも得られる解が一般に近似解となることが知られている。

同等の計算手法は他の分野でも広く用いられている。統計物理のベテ(Bethe)近似(Kabashima and Saad (1998)), また、低密度パリティ検査符号(Gallager (1962), MacKay (1999)), ターボ符号(Berrou et al. (1993), McEliece et al. (1998))といった誤り訂正符号の復号法などは確率伝搬法と等しいことが知られている。統計物理や符号理論で扱う問題に対するグラフィカルモデルは一般にはループを持ち、収束性や近似精度といった問題がある。

確率推論の問題はMCMCなどの手法によっても解が得られる。その場合ループがあっても構わないが、精度の良い解を得るために必要な計算量が多い。前に挙げた誤り訂正符号ではグラフの確率変数の数が通常数百から千程度であり、実時間で復号するには計算量の少ない確率

¹ 統計数理研究所：〒106-8569 東京都港区南麻布 4-6-7

² 東京都立大学大学院 工学研究科：〒192-0397 東京都八王子市南大沢 1-1

³ 理化学研究所 脳科学総合研究センター：〒351-0198 埼玉県和光市広沢 2-1

伝搬法が適している。

これまで我々は情報幾何学(甘利・長岡(1993), Amari and Nagaoka(2000))に基づき確率伝搬法を表現し, 解の安定性, 近似精度といった問題を扱ってきた(Ikeda et al.(2002), 池田 他(2002), Tanaka et al.(2002), Ikeda et al.(2003, 2004a, 2004b))。本稿ではこれまでの結果をまとめ, 新たに得られた 3 次の摂動展開に基づく近似精度の評価について示す。

2. 確率推論の問題と情報幾何的枠組み

2.1 問題の表現

$x = (x_1, \dots, x_n)^T$ を観測できない確率変数, $y = (y_1, \dots, y_m)^T$ を観測された確率変数とする。本稿では簡単のため x_i が 2 値変数, 特に $x_i \in \{-1, +1\}$ の場合を考える。多値変数への拡張は簡単であり, 連続値への拡張も場合によっては可能である(Ikeda et al.(2003, 2004b))。

確率推論の問題は y の条件付きでの x の分布 $q(x|y)$ から(簡単のため以下では $q(x|y)$ を $q(x)$ とかくことにする), x に関する推論を得ることである。1 つの方法は $q(x)$ を最大にする x (MAP 推論: maximum a posteriori)を用いることである。MAP 推論は推論結果が真の x と異なる確率を最小にするが, 探索空間が n とともに指数関数的に増える。ここでは別の推論, 周辺事後確率分布の最大化(MPM 推論: maximization of the posterior marginals)を考える。 $q(x)$ の周辺分布を $q(x_i)$, $i = 1, \dots, n$ とするとき MPM 推論では各成分の推論結果を $q(x_i = +1) \geq q(x_i = -1)$ ならば $\hat{x}_i = +1$, それ以外の場合は $\hat{x}_i = -1$ とする。この結果各 x_i が誤って推論される確率は最小となる。 η_i を $q(x)$ による x_i の期待値とする,

$$\eta_i = E_q[x_i] = \sum_{x_i} x_i q(x_i).$$

MPM 推論は $\hat{x}_i = \text{sgn } \eta_i$ とする。これは周辺分布の積 $\prod_{i=1}^n q(x_i)$ あるいは x の期待値

$$\eta = E_q[x],$$

が分かればただちに計算できる。仮に n とともに指数的に増加しない, 簡単な計算によって η あるいは $\prod_{i=1}^n q(x_i)$ が計算できれば MPM 推論は有効な推論手法となる。本稿で扱う確率伝搬法は簡単な計算によって $\prod_{i=1}^n q(x_i)$ の近似を得ようというものである。

今 x_i は 2 値なので, 全ての $q(x)$ に対し $\ln q(x)$ は $\{x_i\}$ の高々 n 次関数として表現できる。グラフィカルモデルで表現できる問題の多くは確率変数間の相互関係が限られており, 低次の関数で表現できる場合が多い。そこで $\ln q(x)$ を次のように表現する

$$(2.1) \quad \ln q(x) = h \cdot x + \sum_{r=1}^L c_r(x) - \psi_q,$$

ここで $h \cdot x = \sum_i h_i x_i$ は x_i の線形項, $c_r(x)$, $r = 1, \dots, L$, は単項式でも多項式でもよいが高次項を表し, ψ_q は規格化定数の対数である。 $c_r(x)$ はクリーク関数とも呼ばれる。ボルツマンマシンや古典的なスピングラスのモデルでは $c_r(x)$ は x_i の二次関数である。

$$c_r(x) = w_{ij} x_i x_j,$$

ここで r は辺の番号を表すもので, x_i と x_j 間の辺を示す。以下ではこのモデルをボルツマンマシンと呼ぶ(図 1)。

一般のグラフィカルモデル, 特に無向グラフでは $q(x)$ を以下のようにクリーク関数の積として定義することが多い。

$$q(x) = \frac{1}{Z_q} \prod_{i=1}^n \phi_i(x_i) \prod_{r \in \mathcal{C}} \phi_r(x_r),$$

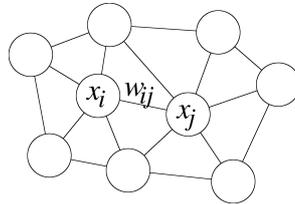


図 1. ボルツマンマシン .

ここで C はクリークの集合である . (2.1) 式の表現と $\phi_i(x_i)$, $\phi_r(\mathbf{x}_r)$ とは次の関係がある .

$$h_i = \frac{1}{2} \ln \frac{\phi_i(x_i = +1)}{\phi_i(x_i = -1)}, \quad c_r(\mathbf{x}) = \ln \phi_r(\mathbf{x}_r), \quad \psi_q = \ln Z_q .$$

ボルツマンマシンでは $\phi_r(\mathbf{x}_r)$ は $\phi_r(x_i, x_j)$ となる . 本論文では $c_r(\mathbf{x})$, $r = 1, \dots, L$ は一次独立であるとする . なお, 本稿では無向グラフのみを扱うが全ての有向グラフは無向グラフで表現できることから一般性は失っていない (Lauritzen and Spiegelhalter (1988)).

2.2 情報幾何学と多様体

本節では情報幾何学に基づく議論のための準備を行う .

まず, 確率伝搬法を考えるために必要な確率分布全体の族 S を考える .

$$(2.2) \quad S = \{p(\mathbf{x}; \theta, \mathbf{v}) \mid p(\mathbf{x}; \theta, \mathbf{v}) = \exp[h \cdot \mathbf{x} + \theta \cdot \mathbf{x} + \mathbf{v} \cdot \mathbf{c}(\mathbf{x}) - \psi(\theta, \mathbf{v})], \theta \in \mathbb{R}^n, \mathbf{v} \in \mathbb{R}^L\},$$

(θ, \mathbf{v}) は自然パラメータであり, $\theta = (\theta_1, \dots, \theta_n)^T$, $\mathbf{v} = (v_1, \dots, v_L)^T$, $\mathbf{c}(\mathbf{x}) = (c_1(\mathbf{x}), \dots, c_L(\mathbf{x}))^T$, $\mathbf{v} \cdot \mathbf{c}(\mathbf{x}) = \sum_{r=1}^L v_r c_r(\mathbf{x})$ である . S は定義より指数型分布族である . $q(\mathbf{x})$ が S に含まれることは $\theta = 0$, $\mathbf{v} = \mathbf{1}_L = (1, \dots, 1)^T$ とおけば簡単に確かめられる . 次に $\mathbf{v} = 0$ となる S の部分多様体を M_0 と呼ぶことにする .

$$M_0 = \{p_0(\mathbf{x}; \theta) = \exp[h \cdot \mathbf{x} + \theta \cdot \mathbf{x} - \psi_0(\theta)] \mid \theta \in \mathbb{R}^n\} .$$

M_0 の分布は各成分が独立であり, 自然パラメータは θ である . さらに, 各成分が独立な分布は全て M_0 に含まれることから $q(\mathbf{x})$ の周辺分布の積 $\prod_{i=1}^n q(x_i)$ も M_0 に含まれる . すなわち $\prod_{i=1}^n q(x_i)$ に対応する M_0 の座標 θ を見つければ MPM 推論が得られる .

次に S 中の e -平坦, m -平坦な部分多様体を定義する .

e -平坦な部分多様体: 部分多様体 $M \subset S$ は次の $r(\mathbf{x}; t)$ が全ての $t \in [0, 1]$, $q(\mathbf{x}), p(\mathbf{x}) \in M$ に対して M に含まれるとき e -平坦である .

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

$c(t)$ は規格化定数である . e -平坦な部分多様体が 1 次元の曲線るとき, 特に e -測地線と呼ぶ . $\{r(\mathbf{x}; t) \mid t \in [0, 1]\}$ は $p(\mathbf{x})$ と $q(\mathbf{x})$ を結ぶ e -測地線である .

m -平坦な部分多様体: 部分多様体 $M \subset S$ は次の $r(\mathbf{x}; t)$ が全ての $t \in [0, 1]$, $q(\mathbf{x}), p(\mathbf{x}) \in M$ に対して M に含まれるとき m -平坦である .

$$r(\mathbf{x}; t) = (1 - t)q(\mathbf{x}) + tp(\mathbf{x}) .$$

m -平坦な部分多様体が 1 次元の曲線るとき, 特に m -測地線と呼ぶ . $\{r(\mathbf{x}; t) \mid t \in [0, 1]\}$ は $p(\mathbf{x})$ と $q(\mathbf{x})$ を結ぶ m -測地線である .

定義よりただちに指数型分布族が e -平坦であることがわかる。したがって S と M_0 は e -平坦である。次に m -射影を定義する。本稿では e -平坦な多様体への m -射影を考えれば十分であるので、次のように定義する。

Definition 1. M を S の e -平坦な部分多様体とし、 $q(x) \in S$ とする。 $q(x)$ から M 上の点 $p(x)$ への m -測地線が Fisher 情報行列を Riemann 計量として M と直交する点を $q(x)$ から M への m -射影と呼ぶ。

Theorem 1. M が S の e -平坦な部分多様体であるとき $q(x) \in S$ から M への m -射影は唯一である。今 $q(x) \in S$ から M への m -射影を $\Pi_M \circ q(x)$ とかくと、 $\Pi_M \circ q(x)$ は M において $q(x)$ からの Kullback-Leibler (KL) divergence を最小にする点である。すなわち

$$(2.3) \quad \Pi_M \circ q(x) = \operatorname{argmin}_{p(x) \in M} D[q(x); p(x)] = \operatorname{argmin}_{p(x) \in M} \sum_x q(x) \ln \frac{q(x)}{p(x)}.$$

ここで $D[\cdot; \cdot]$ は KL-divergence である。

2.3 MPM 推論

本節では MPM 推論が $q(x)$ から M_0 への m -射影と等しいことを示す。(2.3) 式の定義から $q(x)$ から M_0 への m -射影がパラメータ θ^* に対応するならば

$$p_0(x; \theta^*) = \Pi_{M_0} \circ q(x)$$

とかける。これより M_0 への m -射影により求まるパラメータ θ を次のように書くことにする。

$$\theta^* = \pi_{M_0} \circ q(x) = \operatorname{argmin}_{\theta} D[q(x); p_0(x; \theta)].$$

$D[q(x); p_0(x; \theta)]$ を θ で微分し

$$(2.4) \quad \sum_x x q(x) - \partial_{\theta} \psi_0(\theta^*) = 0.$$

∂_{θ} は θ による微分を表す。指数型分布族の定義から

$$(2.5) \quad \partial_{\theta} \psi_0(\theta) = \partial_{\theta} \ln \sum_x \exp(h \cdot x + \theta \cdot x) = \sum_x x p_0(x; \theta).$$

ここで M_0 の期待値パラメータ $\eta_0(\theta)$ を次のように定義する。

$$(2.6) \quad \eta_0(\theta) = \sum_x x p_0(x; \theta) = \partial_{\theta} \psi_0(\theta).$$

M_0 の 2 つの座標系 θ と η_0 とは 1 対 1 に対応するので、(2.4), (2.5), (2.6) 式から m -射影が $q(x)$ の周辺化と等しく、MPM 推論と等価であることがわかる。

3. 確率伝搬法の情報幾何

3.1 確率伝搬法の情報幾何的表現

本節では確率伝搬法の情報幾何的な表現を与える。グラフィカルモデルにおける確率伝搬法はメッセージと呼ばれる変数の更新規則を定義するのが通常である (Pearl (1988), Weiss (2000))。まず、ボルツマンマシンに対する確率伝搬法を Yedidia et al. (2001) にしたがって定義する。確率伝搬法ではメッセージ $m_{ji}(x_i)$ を次の式で更新し、全てが収束した後、得られたメッセージ

$m_{j_i}^*(x_i)$ を用いてビリーフ $b_i(x_i)$ を求める .

$$(3.1) \quad \begin{aligned} m_{j_i}^{t+1}(x_i) &= \frac{1}{Z} \sum_{x_j} \phi_j(x_j) \phi_{ij}(x_i, x_j) \prod_{k \in \mathcal{N}(j) \setminus i} m_{k_j}^t(x_j) \\ b_i(x_i) &= \frac{1}{Z'} \pi_i(x_i) \sum_{k \in \mathcal{N}(i)} m_{k_i}^*(x_i). \end{aligned}$$

$\mathcal{N}(j)$ は j 番目の節につながっている節の集合, Z, Z' はそれぞれ $m_{j_i}^{t+1}(x_i)$ および $b_i(x_i)$ の x_i に関する和を 1 に規格化する . 木のグラフではビリーフ $b_i(x_i)$ は $q(x)$ の x_i に関する周辺分布 $q(x_i)$ と一致するが, ループのある場合には一般に一致せず, 近似となる .

以下ではこの確率伝搬法の情報幾何学に基づく表現を与える . 確率伝搬法では $q(x)$ と $p_0(x; \theta)$ さらに次式の $p_r(x; \zeta_r), r = 1, \dots, L$, を用いる .

$$p_r(x; \zeta_r) = \exp[h \cdot x + c_r(x) + \zeta_r \cdot x - \psi_r(\zeta_r)], \quad \zeta_r \in \mathbb{R}^n, \quad r = 1, \dots, L.$$

$p_r(x; \zeta_r)$ はクリーク関数を 1 つだけ含んだ確率分布である . $c_r(x)$ は一般に $\{x_i\}$ の高次の関数であるが, $\{c_r(x)\}$ の全てを含む場合に比べて 1 つのみを含む場合は扱い易い . 他のクリーク関数は $\zeta_r \cdot x$ で代替されている . $p_r(x; \zeta_r)$ は指数型分布族であり, 以下で定義される M_r は e -平坦である .

$$M_r = \{p_r(x; \zeta_r) \mid \zeta_r \in \mathbb{R}^n\}, \quad r = 1, \dots, L.$$

自然パラメータは ζ_r である . M_r の期待値パラメータを $\eta_r(\zeta_r)$ と定める .

$$(3.2) \quad \eta_r(\zeta_r) = \partial_{\zeta_r} \psi_r(\zeta_r) = \sum_x x p_r(x; \zeta_r), \quad r = 1, \dots, L.$$

$\eta_r(\zeta_r)$ を求めるための計算量は $p_r(x; \zeta_r)$ の周辺化と同等である . このためには $c_r(x)$ に含まれる x_i の全ての組み合わせについて計算しなければいけないが, この計算は可能だとする . 全ての r に対して $\eta_r(\zeta_r)$ を求めるほうが $q(x)$ の周辺化を直接行うより圧倒的に計算量は少ない . なお, ζ_r から $\eta_r(\zeta_r)$ を求めるのは簡単だが, 一般に逆は簡単でないことを注意しておく .

確率伝搬法では $q(x)$ の代わりに $p_r(x; \zeta_r)$ の周辺化を全ての辺の数だけ行い, 全体として $q(x)$ の近似を行う . つまり $p_r(x; \zeta_r), r = 1, \dots, L$ によって各 $c_r(x)$ の影響を表現し, 全ての影響を θ にまとめ, $p_0(x; \theta)$ によって $\prod_i q(x_i)$ を近似する . その際, 繰り返し計算によって $\{\zeta_r\}$ および θ を更新する .

我々は, 誤り訂正符号で用いられる特殊なグラフィカルモデルに対しては池田 他(2002)および Ikeda et al. (2004a)で, 一般のグラフに対しては Ikeda et al. (2003, 2004b)で確率伝搬法の情報幾何的な表現を得ている(Richardson(2000)も同様の枠組みを与えている) . ここではその結果を示す .

$p_0(x; \theta^t)$ および $p_r(x; \zeta_r^t)$ を時刻 t におけるそれぞれ $M_0, M_r, r = 1, \dots, L$ での $q(x)$ の近似とする . 確率伝搬法は ζ_r^t および θ^t を以下のように更新する .

確率伝搬法の情報幾何的表現

- (1) 初期値を $t = 0, \xi_r^t = 0, \zeta_r^t = 0, r = 1, \dots, L$ とする .
- (2) t を 1 つずつ増加させ $\xi_r^{t+1}, r = 1, \dots, L$ を次のように更新する .

$$(3.3) \quad \xi_r^{t+1} = \pi_{M_0} \circ p_r(x; \zeta_r^t) - \zeta_r^t.$$

- (3) θ^{t+1} と ζ_r^{t+1} を以下のように更新する .

$$(3.4) \quad \zeta_r^{t+1} = \sum_{r' \neq r} \xi_{r'}^{t+1}, \quad \theta^{t+1} = \sum_r \xi_r^{t+1} = \frac{1}{L-1} \sum_r \zeta_r^{t+1}.$$

(4) 2 と 3 を $\{\xi_r^t\}$ が収束するまで繰り返す.

$\theta^t = \sum_r \xi_r^t$, $\theta^t = \xi_r^t + \zeta_r^t$ は常に成り立つ. 確率伝搬法が収束した時点でパラメータを θ^* , $\{\zeta_r^*\}$, $\{\xi_r^*\}$ とする. 直感的に確率伝搬法を理解するには $q(x)$, $p_0(x; \theta^*)$, $p_r(x; \zeta_r^*)$, $r = 1, \dots, L$ に関する以下の関係を見るのが良い.

$$\begin{aligned} q(x) &= \exp[\mathbf{h} \cdot \mathbf{x} + c_1(x) + \dots + c_r(x) + \dots + c_L(x) - \psi_q] \\ p_0(x; \theta^*) &= \exp[\mathbf{h} \cdot \mathbf{x} + \xi_1^* \cdot \mathbf{x} + \dots + \xi_r^* \cdot \mathbf{x} + \dots + \xi_L^* \cdot \mathbf{x} - \psi_0(\theta^*)] \\ p_1(x; \zeta_1^*) &= \exp[\mathbf{h} \cdot \mathbf{x} + c_1(x) + \dots + \xi_r^* \cdot \mathbf{x} + \dots + \xi_L^* \cdot \mathbf{x} - \psi_1(\zeta_1^*)] \\ &\vdots \\ p_r(x; \zeta_r^*) &= \exp[\mathbf{h} \cdot \mathbf{x} + \xi_1^* \cdot \mathbf{x} + \dots + c_r(x) + \dots + \xi_L^* \cdot \mathbf{x} - \psi_r(\zeta_r^*)] \\ &\vdots \\ p_L(x; \zeta_L^*) &= \exp[\mathbf{h} \cdot \mathbf{x} + \xi_1^* \cdot \mathbf{x} + \dots + \xi_r^* \cdot \mathbf{x} + \dots + c_L(x) - \psi_L(\zeta_L^*)]. \end{aligned}$$

この式から確率伝搬法では各 $c_r(x)$ を $\xi_r^* \cdot \mathbf{x}$ によって表現していることがわかる.

情報幾何的表現で用いた ξ_r と (3.1) 式で用いたメッセージとの関係をボルツマンマシンの場合に示す. 辺 r が i と j の節を結ぶものとする,

$$\xi_{r,i} = \frac{1}{2} \ln \frac{m_{ji}(x_i = +1)}{m_{ji}(x_i = -1)}, \quad \xi_{r,j} = \frac{1}{2} \ln \frac{m_{ij}(x_j = +1)}{m_{ij}(x_j = -1)}, \quad \xi_{r,k} = 0 \text{ for } k \neq i, j.$$

この関係を用いると (3.1) 式のメッセージの更新則と (3.3) 式が対応すること, $p_0(x; \theta^*)$ が $\prod_i b_i(x_i)$ と等しいことが分る.

3.2 停留点の持つ性質

Theorem 2. 確率伝搬法の停留点は以下の 2 つの条件を満たす (Ikeda et al. (2004a)).

m -条件: $\theta^* = \pi_{M_0} \circ p_r(x; \zeta_r^*)$.

e -条件: $\theta^* = \frac{1}{L-1} \sum_{r=1}^L \zeta_r^*$.

確率伝搬法の停留点で m -条件が満たされていることは (3.3) 式と $\theta^* = \zeta_r^* + \xi_r^*$ から簡単に確かめられる. e -条件は (3.4) 式からただちに確かめられる.

情報幾何的な意味を明らかにするため, S の 2 つの部分多様体 M^* と E^* を定義する.

$$(3.5) \quad \begin{aligned} M^* &= \left\{ p(x) \left| p(x) \in S, \sum_x x p(x) = \sum_x x p_0(x; \theta^*) = \eta_0(\theta^*) \right. \right\}, \\ E^* &= \left\{ p(x) = C p_0(x; \theta^*)^{t_0} \prod_{r=1}^L p_r(x; \zeta_r^*)^{t_r} \left| \sum_{r=0}^L t_r = 1, t_r \in \mathfrak{R} \right. \right\}, \end{aligned}$$

C : 規格化定数.

M^* は m -平坦な部分多様体であり M_0 および M_r , $r = 1, \dots, L$ と直交する. 一方, E^* は e -平坦な部分多様体である. また定義より, M^* は $p_0(x; \theta^*)$ を, E^* は $p_0(x; \theta_0^*)$ と $p_r(x; \zeta_r^*)$, $r = 1, \dots, L$ を含む.

前の 2 つの条件はこれらの部分多様体を用いて次のようにかける.

m -条件: M^* が $p_r(x; \zeta_r^*)$, $r = 1, \dots, L$ を含む.

e -条件: E^* が $q(x)$ を含む .

m -条件が上のように書き直せることは定義より明らかである . e -条件については (3.5) 式で $t_0 = -(L-1), t_1 = \dots = t_L = 1$ とおけば $\theta^* = \sum_{r=1}^L \zeta_r^*/(L-1)$ と同値なことがわかる .

確率伝搬法が収束した点では e -条件と m -条件が同時に満たされるが , だからといって $p_0(x; \theta^*)$ が真の周辺分布 $\prod_{i=1}^n q(x_i)$ であるわけではない . これは M^* と E^* との間の差から生じる . M^* と E^* はともに $p_0(x; \theta^*)$ と $p_r(x; \zeta_r^*)$ を含んでいるが , 他の点については明らかではない . 本来ならば $q(x)$ が M^* に含まれればよいのだが , 計算量を減らすためにこれを E^* で置き換えるのが確率伝搬法である . 同様の仕組みは平均場近似などでも用いられている .

特殊なのはグラフが木の構造の場合である . 木のグラフでは確率伝搬法によって必ず収束し , 正しい周辺分布が求まることが分かっている (Pearl (1988)) . これより次の結果が得られる .

Proposition 1. $q(x)$ が木のグラフで表現できるとき $q(x), p_0(x; \theta^*), p_r(x; \zeta_r^*), r = 1, \dots, L$ が M^* と E^* に含まれる .

木のグラフの場合には $p_0(x; \theta^*) = \prod_i q(x_i)$ であり , 確率伝搬法によって真の周辺分布が求まるが , ループのあるグラフの場合には一般には $q(x) \notin M^*$ である . すなわち確率伝搬法の収束点は特殊な場合を除き , 真の周辺分布を与えない .

3.3 停留点の安定性

ここでは線形近似に基づき , 確率伝搬法の停留点の局所的な安定性を考える . 停留点 $\{\zeta_r^*\}$ に摂動 $\Delta\zeta_r$ を加え $\zeta_r = \zeta_r^* + \Delta\zeta_r$ とする . 確率伝搬法のステップ (2) の結果を $\xi_r = \xi_r^* + \Delta\xi_r$ とすると $\Delta\xi_r$ は次のように表される .

$$\Delta\xi_r = I_0(\theta^*)^{-1} I_r(\zeta_r^*) \Delta\zeta_r - \Delta\zeta_r = (I_0(\theta^*)^{-1} I_r(\zeta_r^*) - E_n) \Delta\zeta_r .$$

ここで $I_0(\theta^*)$ は $p_0(x; \theta)$ の , $I_r(\zeta_r^*)$ は $p_r(x; \zeta_r)$ の Fisher 情報量行列であり E_n は n 次の単位行列である . 確率伝搬法の計算を一ステップ行った後のパラメータを $\zeta_r' = \zeta_r^* + \Delta\zeta_r'$ とかくと ,

$$\Delta\zeta_r' = \sum_{r' \neq r}^L (I_0(\theta^*)^{-1} I_{r'}(\zeta_{r'}^*) - E_n) \Delta\zeta_{r'} .$$

この結果から次の条件が成り立つとき , 確率伝搬法の停留点が安定であることがわかる .

Theorem 3. (Ikeda et al. (2004a)) 停留点の周りでの確率伝搬法の性質は次のように近似できる .

$$\begin{pmatrix} \Delta\zeta_1' \\ \vdots \\ \Delta\zeta_L' \end{pmatrix} = T \begin{pmatrix} \Delta\zeta_1 \\ \vdots \\ \Delta\zeta_L \end{pmatrix} ,$$

ここで

$$T = \begin{pmatrix} O & I_0^{-1} I_2 - E_n & \dots & I_0^{-1} I_L - E_n \\ I_0^{-1} I_1 - E_n & O & & \vdots \\ \vdots & & \ddots & \vdots \\ I_0^{-1} I_1 - E_n & \dots & \dots & O \end{pmatrix} ,$$

である . ただし $I_0 = I_0(\theta^*), I_r = I_r(\zeta_r^*)$ とした . T の全ての固有値 $\lambda_i, i = 1, \dots, nL$ が $|\lambda_i| < 1$ を満たすとき , 停留点は安定である .

上の結果は全てのパラメータを同時に更新する場合の停留点の条件である．実際にはパラメータを一つずつ更新する場合やランダムに更新するなど様々な方法があり，それぞれで安定性は多少異なる．いずれの場合も同様な解析が可能である．

3.4 解の近似精度

本節では摂動展開に基づく近似精度の評価，すなわち確率伝搬法で求めた結果と真の周辺分布との差について述べる．2 次の摂動展開の結果については Ikeda et al. (2004a)にあるが，まずその結果を一般のグラフィカルモデルに対して導き，その後ボルツマンマシンに対する 3 次の摂動展開の結果について述べる．

以下では (2.2) 式の $p(x; \theta, v)$ を用いて議論する．定義より明らかに

$$p_0(x; \theta) = p(x; \theta, \mathbf{0}), \quad p_r(x; \zeta_r) = p(x; \zeta_r, e_r), \quad q(x) = p(x; \mathbf{0}, \mathbf{1}_L)$$

である．ただし

$$e_r = (0, \dots, 0, \underset{\uparrow}{1}, 0, \dots, 0)^T, \quad \mathbf{1}_L = (\underbrace{1, \dots, 1}_L)^T = \sum_{r=1}^L e_r.$$

さらに $p(x; \theta, v)$ の期待値パラメータを次のように定義する

$$\eta(\theta, v) = \partial_\theta \psi(\theta, v) = \sum_x x p(x; \theta, v).$$

近似精度の解析のため $p(x; \theta, v)$ を M^* に拘束する．すなわち

$$\eta(\theta, v) = \eta(\theta^*, \mathbf{0}) = \eta(\theta^*)$$

が常に満たされるようにする．このとき θ は v の関数である．必要があれば $\theta(v)$ とかく．以下では v が $\mathbf{0}$ から $\mathbf{1}_L$ まで変化するにつれて θ が θ^* からどのように変化するかを評価する．まず $\theta(v)$ の 2 次までの摂動展開の結果を示す． M^* に拘束されていることから

$$(3.6) \quad \mathbf{0} = \frac{d}{dv} \eta(\theta, v) = \frac{\partial \eta}{\partial \theta} \frac{\partial \theta}{\partial v} + \frac{\partial \eta}{\partial v}.$$

$\partial \eta / \partial \theta$ および $\partial \eta / \partial v$ は $p(x; \theta, v)$ の Fisher 情報量行列の一部である．これ以降 i, j, k は θ の r, s, t は v の成分を表すことにする． $G_{\theta\theta} = (\partial \eta / \partial \theta)$, $G_{\theta v} = (\partial \eta / \partial v)$ と定める．ここで任意の v と $\theta = \theta(v)$ に対して $G_{\theta\theta} = I_0(\theta)$ であることに注意する．(3.6) 式から

$$(3.7) \quad \mathbf{0} = I_0(\theta) \frac{\partial \theta}{\partial v} + G_{\theta v}(\theta), \quad \frac{\partial \theta}{\partial v} = -I_0^{-1}(\theta) G_{\theta v}(\theta).$$

これにより $\theta(v)$ の v に関する 1 次の微分 $\partial \theta / \partial v$ が得られた．また $\tilde{G}_{\theta v}$ を $-\partial \theta / \partial v$ と定義する．2 次の微分係数は

$$\frac{d^2}{dv dv} \eta(\theta, v) = 0,$$

より

$$(3.8) \quad I_0(\theta) \frac{\partial^2 \theta}{\partial v \partial v'} = -T_{\theta v v'} - T_{\theta \theta \theta} \frac{\partial \theta}{\partial v} \frac{\partial \theta}{\partial v'} - T_{\theta \theta v} \frac{\partial \theta}{\partial v'} - T_{\theta \theta v'} \frac{\partial \theta}{\partial v}$$

となる．ここで

$$T_{\theta \theta \theta} = \frac{\partial^3 \psi}{\partial \theta \partial \theta \partial \theta}, \quad T_{\theta \theta v} = \frac{\partial^3 \psi}{\partial \theta \partial \theta \partial v}, \quad T_{\theta v v'} = \frac{\partial^3 \psi}{\partial \theta \partial v \partial v'}$$

である. $\partial\theta/\partial v = -\tilde{G}_{\theta v}(\theta)$ および $\partial^2\theta/\partial v^2$ を $(\theta, v) = (\theta^*, 0)$ で評価し, $\theta(v)$ をまずは 2 次のテーラー展開により v で近似する. $(\theta, v) = (\theta^*, 0)$ におけるオペレータ d/dv を以下のように定義すれば

$$\frac{d}{dv} = B = \frac{\partial}{\partial v} - \tilde{G}_{\theta v}(\theta^*) \frac{\partial}{\partial \theta}.$$

$\partial^2\theta/\partial v\partial v$ は

$$\left. \frac{\partial^2\theta}{\partial v\partial v} \right|_{v=0} = -I_0(\theta^*)^{-1} B^2 \eta(\theta^*),$$

となる. 簡単のため B^2 の (r, s) 成分を $B_{rs} = B_r B_s$ とかくことにする. なお $(d^2/dv dv)\eta(\theta^*) = 0$ であっても一般に $B^2\eta(\theta^*) \neq 0$ である.

$\theta(v)$ を $(\theta^*, 0)$ の周りで v の 2 次までで近似した結果は以下の通りである.

$$(3.9) \quad \begin{aligned} \theta(v) &= \theta^* + \left. \frac{\partial\theta}{\partial v} \right|_{v=0} v + \frac{1}{2} v^T \left. \frac{\partial^2\theta}{\partial v\partial v} \right|_{v=0} v + o(\|v\|^3) \\ &\simeq \theta^* - \tilde{G}_{\theta v}(\theta^*) v - \frac{1}{2} v^T I_0^{-1}(\theta^*) (B^2 \eta(\theta^*)) v. \end{aligned}$$

ここまで m -条件を考えていたが, e -条件を考慮する必要がある. e -条件は

$$(3.10) \quad \theta^* = - \sum_{r=1}^L (\zeta_r^* - \theta^*),$$

とかける. 次の分布を考える

$$(3.11) \quad p(x; \zeta_r, \epsilon e_r) = \exp[h \cdot x + \zeta_r \cdot x + \epsilon c_r(x) - \psi(\zeta_r, \epsilon e_r)].$$

ここで $p(x; \zeta_r, \epsilon e_r)|_{\epsilon=1} = p_r(x; \zeta_r)$ である. $p(x; \zeta_r, \epsilon e_r)$, $r = 1, \dots, L$ が M^* に含まれるように拘束すると (3.9) 式の結果から $\zeta_r - \theta^*$ は ϵ によって以下のように近似できる.

$$\zeta_r - \theta^* \simeq -\tilde{G}_{\theta v}(\theta^*) e_r \epsilon - \frac{1}{2} I_0^{-1}(\theta^*) B_{rr} \eta(\theta^*) \epsilon^2.$$

以下の議論では θ に対する $c_r(x)$ の影響は小さいと仮定する. この仮定の下では ϵ がある程度大きな値を取っても近似が成り立つことから $\epsilon \rightarrow 1$ とした点を考える. この意味での近似を \approx によって表すこととする. $\epsilon \rightarrow 1$ ととれば $\zeta_r - \theta^*$ は $\zeta_r^* - \theta^*$ となり

$$\zeta_r^* - \theta^* \approx -\tilde{G}_{\theta v}(\theta^*) e_r - \frac{1}{2} I_0^{-1}(\theta^*) B_{rr} \eta(\theta^*).$$

(3.10) 式から θ^* は次式を満たす.

$$(3.12) \quad \theta^* = - \sum_{r=1}^L (\zeta_r^* - \theta^*) \approx \tilde{G}_{\theta v}(\theta^*) \mathbf{1}_L + \frac{1}{2} I_0^{-1}(\theta^*) \sum_r B_{rr} \eta(\theta^*).$$

ここで次の分布を考える.

$$(3.13) \quad p(x; u, \epsilon \mathbf{1}_L) = \exp[h \cdot x + u \cdot x + \epsilon \mathbf{1}_L \cdot c(x) - \psi(u, \epsilon \mathbf{1}_L)].$$

この分布を M^* に拘束し ϵ を 0 から 1 まで増やすと u は $\theta(\mathbf{1}_L)$ となる. 一般に $\theta(\mathbf{1}_L) \neq 0$ であることから $p(x; \theta(\mathbf{1}_L), \mathbf{1}_L)$ は $q(x)$ と等しくない. (3.9) 式の結果から

$$\theta(\mathbf{1}_L) - \theta^* \approx -\tilde{G}_{\theta v}(\theta^*) \mathbf{1}_L - \frac{1}{2} I_0^{-1}(\theta^*) \sum_{r,s} B_{rs} \eta(\theta^*).$$

(3.12) 式から

$$(3.14) \quad \theta(\mathbf{1}_L) \approx -\frac{1}{2} I_0^{-1}(\theta^*) \sum_{r \neq s} B_{rs} \eta(\theta^*).$$

期待値パラメータはテーラー展開に基づき

$$(3.15) \quad \eta(0, \mathbf{1}_L) \simeq \eta(\theta(\mathbf{1}_L), \mathbf{1}_L) - \nabla_{\theta} \eta(\theta^*) \theta(\mathbf{1}_L) \approx \eta(\theta^*) + \frac{1}{2} \sum_{r \neq s} B_{rs} \eta(\theta^*)$$

と近似できる。 $\eta(0, \mathbf{1}_L)$ は $q(x)$ による x の期待値である。したがって (3.15) 式は $q(x)$ による x の期待値と確率伝搬法による結果との差を示している。

Theorem 4. 真の条件付確率 $q(x)$ による x の期待値を $\eta_{MPM} = \eta(0, \mathbf{1}_L)$ とし、確率伝搬法によって求めた期待値を $\eta(\theta^*)$ とする。2 次までの摂動展開の結果 η_{MPM} と $\eta(\theta^*)$ の差は以下のように近似できる。

$$(3.16) \quad \eta_{MPM} \approx \eta(\theta^*) + \frac{1}{2} \sum_{r \neq s} B_{rs} \eta(\theta^*).$$

この定理に示された近似誤差は M^* の e -曲率と関係している (Ikeda et al. (2004a))。グラフィカルモデルが与えられればこの項を具体的に計算できる。具体的な形はモデルによって異なるが、 $c_r(x)$ が $\{x_i\}$ の単項式の場合には次の定理が得られている (Tanaka et al. (2002))。

Theorem 5. $c_r(x)$, $c_{r'}(x)$, $r \neq r'$ が共通な x_i を 2 つ以上含まないとき、摂動展開による 2 次の誤差項 $\sum_{r \neq s} B_{rs} \eta(\theta^*)/2$ は 0 となる。

ボルツマンマシンでは $c_r(x) = w_{ij} x_i x_j$ となるため、 $c_r(x)$, $c_{r'}(x)$, $r \neq r'$ は共通の x_i を持っていたとしても一つのみである。したがってこの近似誤差項は 0 となることがわかる。

ボルツマンマシンの誤差項を評価するためには少なくとも 3 次の摂動展開を行う必要がある。ここでは最近得た結果を示す。3 次までの近似は

$$(3.17) \quad \theta(v) \simeq \theta^* + \frac{\partial \theta}{\partial v} \Big|_{v=0} v + \frac{1}{2} v^T \frac{\partial^2 \theta}{\partial v \partial v} \Big|_{v=0} v + \frac{1}{6} \sum_{r,s,t} \frac{\partial^3 \theta}{\partial v_r \partial v_s \partial v_t} \Big|_{v=0} v_r v_s v_t$$

から求まるわけだが、(3.14) 式の誤差項については 2 次までの項が 0 となることから

$$\theta(\mathbf{1}_L) \simeq \frac{1}{6} \left[\sum_{r,s,t} \frac{\partial^3 \theta}{\partial v_r \partial v_s \partial v_t} \Big|_{v=0} - \sum_r \frac{\partial^3 \theta}{\partial v_r^3} \Big|_{v=0} \right]$$

となる。 $\partial^3 \theta / \partial v_r \partial v_s \partial v_t$ は

$$\frac{d^3}{dv dv dv} \eta(\theta, v) = 0,$$

から多少の計算をすれば求まる。ボルツマンマシンに対して具体的に計算をすると、次の結果を得る。

$$\theta_i(\mathbf{1}_L) \approx \eta_i^* \sum_{j \neq k} w_{ij} w_{jk} w_{ki} (1 - \eta_j^{*2}) (1 - \eta_k^{*2}).$$

これは一番短いループのみが $\theta_i(\mathbf{1}_L)$ に寄与することを示している (図 2)。

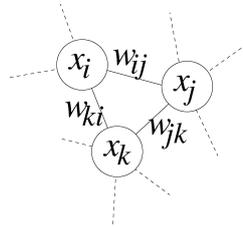


図 2. 誤差項へ寄与する短いループ.

Theorem 6. 3 次までの摂動展開により, ボルツマンマシンに対する η_{MPM} と $\eta(\theta^*)$ の差は以下のように近似できる.

$$(3.18) \quad \eta_{MPM,i} \approx \eta_i^* - \eta_i^*(1 - \eta_i^{*2}) \sum_{j \neq k} w_{ij} w_{jk} w_{ki} (1 - \eta_j^{*2})(1 - \eta_k^{*2}).$$

現在は 3 次までの項の結果しかないが, さらに高次の摂動展開はより長いループを反映した項が現われるものと考えられる.

(3.16), (3.18) 式の結果を用いると MPM 解の精度が向上できる期待がある. しかし, 例えばボルツマンマシンの場合に計算してみると, 必ずしも精度が向上するとも限らない. 原因としてはいくつか考えられる. 摂動展開の仮定としている θ への $c_r(x)$ の影響が小さいとするが, w_{ij} がある程度大きい場合にはこの近似が成り立たない, また, ボルツマンマシンなどではより高次の項の影響があり, 3 次の摂動展開の項のみを考えることはあまり意味がない, などである. 精度は確率伝搬法の一つの重要な問題である. 今後, どうやって 2 次, 3 次, あるいはより高次の摂動展開と精度向上を結びつけるかが課題の一つである.

4. まとめ

確率伝搬法は統計物理学や誤り訂正符号など, 応用上重要な手法である. 本稿では確率伝搬法を理解するための情報幾何学的な枠組みについて述べた.

一般に確率伝搬法で重要な問題は収束性と近似精度である. 我々は提案した枠組みに基づき, 局所的な安定性の条件を示し, 近似誤差についても摂動展開から主要項を示した. 近似誤差については 3 次までの漸近展開の結果では最も短いループが主要項を構成しているが, 今後より高次の展開についても調べる必要がある.

確率伝搬法は様々な分野で独自の研究がなされており, 関連するアルゴリズムや近似についても広く論じられている. 情報幾何学による枠組みではそれらを全て表現し, 等しく扱うことが可能である. 我々は CCCP, TRP, GBP といった関連するアルゴリズムがこの枠組みで表現できることを示した(Ikeda et al. (2004b)). 今後は各分野で得られた結果を理解し, 新たなアルゴリズムの提案に結び付けたいと考えている.

謝 辞

査読者には有益な御指摘を頂きました. 感謝致します.

参 考 文 献

- 甘利俊一, 長岡浩司(1993). 『情報幾何の方法』, 岩波講座 応用数学 [対象 12], 岩波書店, 東京.
- Amari, S. and Nagaoka, H. (2000). *Methods of Information Geometry*, AMS and Oxford University Press, Providence, Rhode Island.
- Berrou, C., Glavieux, A. and Thitimajshima, P. (1993). Near Shannon limit error-correcting coding and decoding: Turbo-codes, *Proceedings of IEEE International Conference on Communications*, Geneva, Switzerland.
- Gallager, R. G. (1962). Low density parity check codes, *IRE Transactions on Information Theory*, **IT-8**, 21–28.
- 池田思朗, 田中利幸, 甘利俊一(2002). ターボ復号の情報幾何, 電子情報通信学会論文誌, **J85-D-II(5)**, 758–765.
- Ikeda, S., Tanaka, T. and Amari, S. (2002). Information geometrical framework for analyzing belief propagation decoder, *Advances in Neural Information Processing Systems*, Vol. 14 (eds. T. G. Dietterich, S. Becker and Z. Ghahramani), 407–414, MIT Press, Cambridge, Massachusetts.
- Ikeda, S., Tanaka, T. and Amari, S. (2003). Stochastic reasoning, free energy and information geometry, Research Memorandum, No. 890, The Institute of Statistical Mathematics, Tokyo.
- Ikeda, S., Tanaka, T. and Amari, S. (2004a). Information geometry of turbo codes and low-density parity-check codes, *IEEE Transactions on Information Theory*, **50(6)**, 1097–1114.
- Ikeda, S., Tanaka, T. and Amari, S. (2004b). Stochastic reasoning, free energy and information geometry, *Neural Computation*, **16(9)**, 1779–1810.
- Jordan, M. I. (1999). *Learning in Graphical Models*, MIT Press, Cambridge, Massachusetts.
- Kabashima, Y. and Saad, D. (1998). Belief propagation vs. TAP for decoding corrupted messages, *Europhysics Letters*, **44(5)**, 668–674.
- Lauritzen, S. L. and Spiegelhalter, D. J. (1988). Local computations with probabilities on graphical structures and their application to expert systems, *Journal of the Royal Statistical Society B*, **50**, 157–224.
- MacKay, D. J. C. (1999). Good error-correcting codes based on very sparse matrices, *IEEE Transactions on Information Theory*, **45(2)**, 399–431.
- McEliece, R. J., MacKay, D. J. C. and Cheng, J.-F. (1998). Turbo decoding as an instance of Pearl's "belief propagation" algorithm, *IEEE Journal on Selected Areas in Communications*, **16(2)**, 140–152.
- Pearl, J. (1988). *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, Morgan Kaufmann, San Mateo, California.
- Richardson, T. J. (2000). The geometry of turbo-decoding dynamics, *IEEE Transactions on Information Theory*, **46(1)**, 9–23.
- Tanaka, T., Ikeda, S. and Amari, S. (2002). Information-geometrical significance of sparsity in Gallager codes, *Advances in Neural Information Processing Systems*, Vol. 14 (eds. T. G. Dietterich, S. Becker and Z. Ghahramani), 527–534, MIT Press, Cambridge, Massachusetts.
- Weiss, Y. (2000). Correctness of local probability propagation in graphical models with loops, *Neural Computation*, **12(1)**, 1–41.
- Yedidia, J. S., Freeman, W. T. and Weiss, Y. (2001). Bethe free energy, Kikuchi approximation, and belief propagation algorithms, Technical Report 2001-16, Mitsubishi Electric Research Laboratories, Cambridge, Massachusetts.

Information Geometrical Framework to Analyze Belief Propagation Algorithm

Shiro Ikeda¹, Toshiyuki Tanaka² and Shun-ichi Amari³

¹The Institute of Statistical Mathematics

²Department of Electronics and Information Engineering, Tokyo Metropolitan University

³Brain Science Institute, RIKEN

Belief propagation (BP) is a universal method of stochastic reasoning. It gives exact inference for stochastic models with tree interactions, and works well even if the models have loopy interactions. Its performance has been analyzed separately in many fields, such as, AI, statistical physics, information theory, and information geometry. The present paper provides a unified framework for understanding BP. The stability of BP is analyzed from this framework, and its approximation accuracy is investigated.