

オンライン学習における適応自然勾配法の評価

井上 真郷^{†,††} 朴 慧暎[†] 岡田 真人^{†,†††}

† 理化学研究所 脳科学総合研究センター 脳数理研究チーム
〒 351-0198 埼玉県和光市広沢 2-1

†† 京都大学大学院医学研究科 耳鼻咽喉科・頭頸部外科学
〒 606-8507 京都府京都市左京区聖護院河原町 54

††† 科学技術振興事業団 戦略的創造研究推進事業 (さきがけ研究 21) 「協調と制御」研究領域
〒 351-0198 埼玉県和光市広沢 2-1

E-mail: †{minoue,hypark,okada}@brain.riken.go.jp

あらまし 適応自然勾配法は自然勾配法における Fisher 情報行列の逆行列を近似的に求めるもので、1) 入力分布を予め知っている必要が無い、2) 入力次元を N とすると、 $O(N^3)$ だった計算コストを $O(N^2)$ にする、という利点がある。本研究ではソフトコミティーマシンでのオンライン学習の枠組みで、統計力学的手法を用いて、適応自然勾配法の性能を評価した。その結果、ネットワークパラメータの更新割合と Fisher 情報行列の逆行列の推定値の更新割合との比を小さくすれば、自然勾配法とほぼ同等の性能が出ることが分かった。

キーワード 適応自然勾配法, 自然勾配法, 多層パーセプトロン, ソフトコミティーマシン, プラトー, 統計力学

Evaluation of Adaptive Natural Gradient Descent in On-Line Learning

Masato INOUE^{†,††}, Hyeyoung PARK[†], and Masato OKADA^{†,†††}

† Laboratory for Mathematical Neuroscience, RIKEN Brain Science Institute, Saitama 351-0198, Japan

†† Department of Otolaryngology, Head and Neck Surgery, Graduate School of Medicine,
Kyoto University, Kyoto 606-8507, Japan

††† “Intelligent Cooperation and Control,” PRESTO, JST, Saitama 351-0198

E-mail: †{minoue,hypark,okada}@brain.riken.go.jp

Abstract Adaptive natural gradient descent (ANGD) approximates the inverse of Fisher information matrix directly. Compared to natural gradient descent (NGD), it has following advantages ; 1) it works without any knowledge of the input distribution, and 2) it reduces the calculation cost from $O(N^3)$ to $O(N^2)$. In this report, we evaluate ANGD using statistical mechanics techniques in respect of on-line learning of soft committee machines. We also show that it has almost the same performance to NGD if the ratio of the update speed of network parameters to the one of the estimator of the inverse of Fisher information matrix is set small enough.

Key words adaptive natural gradient descent, natural gradient descent, multilayer perceptron, soft committee machine, plateau, statistical mechanics

1. はじめに

多層パーセプトロンを、誤差関数を定義して勾配法で学習させる場合、隠れユニットの入れ換え対称性に起因する鞍点構造が生ずる [1], [2]. 特に最急降下法では、学習の過程で鞍点近くを通過するため、この時に学習曲線のプラトーが生じ、大きな問題となっている. 互いに同構造の教師ネットワークと生徒ネットワークを用意し、教師のパラメータを教師が出す出力をもとに生徒が学ぶというケースで議論すると、教師の隠れユニット間で、入力に対する重みベクトルに強い相関がある場合、このプラトーは更に増加することが分かっている [6]. 一方、自然勾配法 [3]~[5] では、学習割合を十分小さく設定すれば、教師の隠れユニット同士の相関に関係なく、プラトーは殆ど生じないという望ましい性質がある [6].

しかし、自然勾配法は反面、1) 特異点付近での学習が不安定になる、2) 入力の分布を予め知っている必要がある、3) Fisher 情報行列の逆行列の計算が時間がかかる、等の欠点知られている. この内、1) は依然未知の部分が多いが、2,3) については、適応自然勾配法が代替方法として提案されている [7], [8]. そこで今回我々は、自然勾配法に見られた望ましい性質が、適応自然勾配法の近似によってどの程度変化しているかを評価した.

様々な学習手法は一般にバッチとオンラインの二通りの学習のさせ方が存在する. オンライン学習では学習サンプルを一度しか使用しないため、系の状態と学習サンプルが独立であるため、解析が容易になるメリットがある. 更に統計力学的手法を用いると、本来確率的にしか表されない学習のダイナミクスが、入力次元 $N \rightarrow \infty$ で漸近的に決定論的になるため、確率的シミュレーションを繰返し行わなくとも、一般的な性質を引き出し易くなる [1].

今回我々は、多層パーセプトロンを簡略化したソフトコミティーマシンのオンライン学習の枠組みで、自然勾配法と適応自然勾配法の違いを、統計力学的手法を用いて評価したので報告する.

2. モデル

教師ネットワークは中間層のユニット数 M 個、生徒ネットワークは K 個のソフトコミティーマシンを考える (図 1).

$$\zeta \equiv f_B(\xi), \quad f_B(\xi) \equiv \sum_{k=1}^M g(B_k^T \xi), \quad (1)$$

$$\zeta' \equiv f_J(\xi) + n, \quad f_J(\xi) \equiv \sum_{k=1}^K g(J_k^T \xi), \quad (2)$$

但し、 \bullet^T は転置を表す. ここで、 $B_i \in \mathbb{R}^N$ 及び $J_i \in \mathbb{R}^N$ は、 i 番目のユニットの、入力 $\xi \in \mathbb{R}^N$ に対する重みベクトルである. 教師ネットワークの確率分布は、入力分布を $p(\xi)$ とすると、

$$p(\xi, \zeta) \equiv p(\xi) \delta(\zeta - f_B(\xi)) \quad (3)$$

となる. 但し、 $\delta(\bullet)$ はデルタ関数である. 生徒ネットワークには、正規分布ノイズ $n \sim \mathcal{N}(0, \sigma^2)$ が加えられることとすると、確率分布は、

$$p_J(\xi, \zeta') \equiv p(\xi) p_J(\zeta' | \xi) \quad (4)$$

$$p_J(\zeta' | \xi) \equiv \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\{\zeta' - f_J(\xi)\}^2}{2\sigma^2}\right) \quad (5)$$

となる. 学習によってパラメータベクトル $J \equiv [J_1^T, J_2^T, \dots, J_K^T]^T$

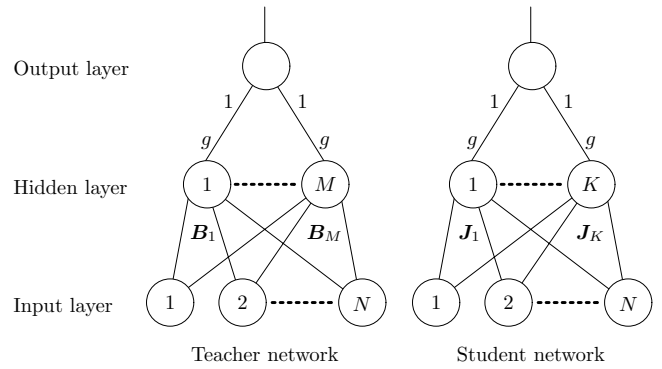


図 1 ソフトコミティーマシン

$\in \mathbb{R}^{KN}$ を少しずつ更新することで、生徒の確率分布を教師の確率分布に近づけることを考える.

与えられた個々の正解サンプル $\{\xi, \zeta\}$ に対する誤差関数を、対数損失を用いて

$$\epsilon_J(\xi, \zeta) \equiv -\ln p_J(\zeta | \xi) + c_0 = \frac{1}{2\sigma^2} \{\epsilon_d(\xi)\}^2 \quad (6)$$

$$\epsilon_d(\xi) \equiv f_J(\xi) - f_B(\xi) \quad (7)$$

で定義する. 但し、 $c_0 \equiv -\ln \sqrt{2\pi\sigma^2}$ はパラメータが完全に一致している ($J = B$) 時の誤差を 0 にする定数である. また、汎化誤差を

$$\epsilon_g(J) \equiv \langle \epsilon_J(\xi, \zeta) \rangle_{\{\xi, \zeta\}} \quad (8)$$

で定義する. 但し、 $\langle \bullet \rangle_{\{o\}}$ は、確率変数 o での平均である. 勾配法のパラメータ更新則は一般に、

$$\Delta J = -\frac{\eta}{N} M \nabla_J \epsilon_J(\xi, \zeta) \quad (9)$$

で表され、 $M = I$ (単位行列) とするのが通常. 最急降下法、 $M = G^{-1}$ (G は J の Fisher 情報行列) が自然勾配法 [3] で、 G^{-1} の代わりにその推定値 \hat{G}^{-1} を用いるのが適応自然勾配法 [7], [8] である. Fisher 情報行列は、

$$G \equiv \langle [\nabla_J \ln p_J(\xi, \zeta')] [\nabla_J \ln p_J(\xi, \zeta')]^T \rangle_{\{\xi, \zeta'\}} \\ = \frac{1}{\sigma^2} \langle [\nabla f] [\nabla f]^T \rangle_{\{\xi\}} \quad (10)$$

で定義される. 但し、 ∇f は、 $\nabla_J f_J(\xi)$ の意である. G をオンラインで推定するには、与えられる個々の ξ に対して、更新割合 ($0 < \frac{\rho}{N} < 1$) を設定して、

$$\Delta \hat{G} \equiv \frac{\rho}{N} \left[-\hat{G} + \frac{1}{\sigma^2} [\nabla f] [\nabla f]^T \right] \quad (11)$$

で、更新していく方法が考えられる. この時、 \hat{G}^{-1} の差分は Sherman-Morrison 公式を用いて正確に、

$$\Delta \hat{G}^{-1} = \frac{\frac{\rho}{N}}{1 - \frac{\rho}{N}} \left[\hat{G}^{-1} - \frac{\frac{1}{\sigma^2} \hat{G}^{-1} [\nabla f] [\nabla f]^T \hat{G}^{-1}}{1 + \frac{\rho}{N} \left\{ \frac{1}{\sigma^2} [\nabla f]^T \hat{G}^{-1} [\nabla f] - 1 \right\}} \right] \quad (12)$$

となる. これが適応自然勾配法で、 G でなく G^{-1} を直接推定できるため、逆行列を計算する手間が省ける (逆行列の計算は $O((NK)^3)$, 適応自然勾配法の計算は $O((NK)^2)$ のコストである). また、入力分布 $p(\xi)$ を必要としない. \hat{G}^{-1} の初期値は正定値対称行列であればよいが、 $\sigma^2 I$ とすると、汎化誤差を除いて系が σ^2 に依存しなくなる.

3. 理 論

3.1 オーダーパラメータ

入力を $\xi \sim \mathcal{N}(0, I)$ とすると、系の N 次元ベクトル全て (ξ, B_i, J_i) を原点周りで回転させても本質的に同等であるので、これらの相関のみで系を記述できる [1], [4] . 即ち、

$$\begin{aligned} x &\equiv [x_i]_{i=1, \dots, K}, & x_i &\equiv J_i^T \xi \\ y &\equiv [y_i]_{i=1, \dots, M}, & y_i &\equiv B_i^T \xi \\ Q &\equiv [Q_{ij}]_{i,j=1, \dots, K}, & Q_{ij} &\equiv J_i^T J_j \\ R &\equiv [R_{ij}]_{i=1, \dots, K, j=1, \dots, M}, & R_{ij} &\equiv J_i^T B_j \\ T &\equiv [T_{ij}]_{i,j=1, \dots, M}, & T_{ij} &\equiv B_i^T B_j \end{aligned} \quad (13)$$

である . $f_J(\xi), f_B(\xi)$ は ξ の代わりに x, y を用いてそれぞれ計算できる . また、 Q は生徒の重みベクトル同士の相関、 T は教師の重みベクトル同士の相関、 R は生徒と教師の重みベクトル間の相関である . これにより、入力次元 N が大きい時でも、系のパラメータ数を N に依存しない形に圧縮できる . また、

$$z \equiv [x^T, y^T]^T, \quad C \equiv \begin{bmatrix} Q & R \\ R^T & T \end{bmatrix} \quad (14)$$

とすると、 $z \sim \mathcal{N}(0, C)$ となる .

汎化誤差は、オーダーパラメータで表すことができる .

$$\epsilon_g = \frac{1}{2\sigma^2} \int p(z) dz \{ \epsilon_d(z) \}^2 \quad (15)$$

$$\epsilon_d(z) = f_J(x) - f_B(y) = \sum_{k=1}^{K+M} a_k g(z_k) \quad (16)$$

$$a_i \equiv \begin{cases} 1 & (i = 1, \dots, K) \\ -1 & (i = K+1, \dots, K+M) \end{cases} \quad (17)$$

また、ここから先 $g(x) \equiv \text{erf}(x/\sqrt{2})$ 、 g の一次導関数を $g'(x) = \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}}$ とすると、汎化誤差は解析的に求まる (付録 1 参照) .

$$\epsilon_g = \frac{1}{\pi\sigma^2} \sum_{i,j=1}^{K+M} a_i a_j \arcsin \frac{C_{ij}}{\sqrt{\{C_{ii}+1\}\{C_{jj}+1\}}} \quad (18)$$

3.2 最急降下法

オーダーパラメータの更新式も、オーダーパラメータで表すことができる [1] .

$$\nabla_{J_i} \epsilon_J = \frac{1}{\sigma^2} \epsilon_d(z) \nabla_{J_i} f_J \quad (19)$$

$$\nabla_{J_i} f_J = g'(x_i) \xi \quad (20)$$

となるので、最急降下法の場合の J の更新則は

$$\Delta J_i = -\frac{\eta}{\sigma^2 N} \epsilon_d(z) g'(x_i) \xi \quad (21)$$

となり、オーダーパラメータの更新式は

$$\Delta R_{ij} = [\Delta J_i^T] B_j = -\frac{\eta}{\sigma^2 N} \epsilon_d(z) g'(x_i) y_j \quad (22)$$

$$\Delta S_{ij} \equiv [\Delta J_i^T] J_j = -\frac{\eta}{\sigma^2 N} \epsilon_d(z) g'(x_i) x_j$$

$$\begin{aligned} \Delta Q_{ij} &= [\Delta J_i^T] J_j + [\Delta J_j^T] J_i + [\Delta J_i^T] [\Delta J_j] \\ &= \Delta S_{ij} + \Delta S_{ji} + \frac{\eta^2 \{ \xi^T \xi \}}{\sigma^4 N^2} \{ \epsilon_d(z) \}^2 g'(x_i) g'(x_j) \end{aligned}$$

となる . (ΔS_{ij} は便宜上導入した . 以下も S については同様 .)

次に、時間 α を導入し、一回の更新で $1/N$ の時間が経過するものとし、 $N \rightarrow \infty$ の極限で、系のパラメータがどのように変化するかを考える . 極限をとることで、系のダイナミクスが一意に決まるので、解析が容易になる . 具体的には、例えば ΔR_{ij} は確率変数であるが、現在時刻 α より微小時間 τ 後の値は、他のパラメータが固定されているものとして、

$$\begin{aligned} R_{ij}(\alpha + \tau) &= R_{ij}(\alpha) + \lim_{N \rightarrow \infty} \sum_{k=0}^{\tau N - 1} \Delta R_{ij} \\ &= R_{ij}(\alpha) + \tau \left\langle -\frac{\eta}{\sigma^2} \epsilon_d(z) g'(x_i) y_j \right\rangle_{\{z\}} \end{aligned} \quad (23)$$

と、確率変数でなくなる . 従って、現在の微分値は、

$$\begin{aligned} \frac{\partial R_{ij}(\alpha)}{\partial \alpha} &= \lim_{\tau \rightarrow 0} \frac{R_{ij}(\alpha + \tau) - R_{ij}(\alpha)}{\tau} \\ &= \left\langle -\frac{\eta}{\sigma^2} \epsilon_d(z) g'(x_i) y_j \right\rangle_{\{z\}} \end{aligned} \quad (24)$$

となる . 他も同様で、簡潔のため行列形式で表すと、

$$\begin{aligned} \frac{\partial R}{\partial \alpha} &= -\frac{\eta}{\sigma^2} \langle \delta y^T \rangle_{\{z\}} \\ \frac{\partial S}{\partial \alpha} &\equiv -\frac{\eta}{\sigma^2} \langle \delta x^T \rangle_{\{z\}} \\ \frac{\partial Q}{\partial \alpha} &= \frac{\partial S}{\partial \alpha} + \frac{\partial S^T}{\partial \alpha} + \frac{\eta^2}{\sigma^4} \langle \delta \delta^T \rangle_{\{z\}} \end{aligned} \quad (25)$$

となる . 但し、

$$\delta \equiv [\delta_i]_{i=1, \dots, K}, \quad \delta_i \equiv \epsilon_d(z) g'(x_i) \quad (26)$$

とした . また、 ΔQ_{ij} に残っていた ξ は、 $N \rightarrow \infty$ で $\xi^T \xi \rightarrow N$ となる . 平均の計算については付録 1 参照 .

3.3 自然勾配法

自然勾配法の場合も同様に計算できる [5] . まず $G = [G_{ij}]_{i,j=1, \dots, K}$ をオーダーパラメータを用いて表す . これは、

$$\begin{aligned} \Lambda &\equiv [\lambda_{ij}]_{i,j=1, \dots, K} \\ \Lambda &\equiv [\Lambda_{ij}]_{i,j=1, \dots, K}, \quad \tilde{\Lambda} \equiv [\tilde{\Lambda}_{ij}]_{i,j=1, \dots, K} \end{aligned} \quad (27)$$

$$\Lambda_{ij} = [\tilde{\Lambda}_{ij} - \lambda_{ij} I] Q^{-1}, \quad \tilde{\Lambda}_{ij} = \Lambda_{ij} Q + \lambda_{ij} I \quad (28)$$

を用いて、

$$G_{ij} = \sigma^2 [\lambda_{ij} I + U \Lambda_{ij} U^T] \quad (29)$$

$$= \sigma^2 \left[\lambda_{ij} [I - UV^T] + U \tilde{\Lambda}_{ij} V^T \right] \quad (30)$$

の二形式で表せる . また $G^{-1} = [G_{ij}^{-1}]_{i,j=1, \dots, K}$ も、

$$\begin{aligned} \Theta &\equiv [\theta_{ij}]_{i,j=1, \dots, K} \\ \Theta &\equiv [\Theta_{ij}]_{i,j=1, \dots, K}, \quad \tilde{\Theta} \equiv [\tilde{\Theta}_{ij}]_{i,j=1, \dots, K} \end{aligned} \quad (31)$$

$$\Theta_{ij} = [\tilde{\Theta}_{ij} - \theta_{ij} I] Q^{-1}, \quad \tilde{\Theta}_{ij} = \Theta_{ij} Q + \theta_{ij} I \quad (32)$$

を用いて、

$$G_{ij}^{-1} = \sigma^2 [\theta_{ij} I + U \Theta_{ij} U^T] \quad (33)$$

$$= \sigma^2 \left[\theta_{ij} [I - UV^T] + U \tilde{\Theta}_{ij} V^T \right] \quad (34)$$

の二形式で表せることが分かっている [4] . 但し、 $U \equiv$

$[J_1 \cdots J_K]$ は $N \times K$ 行列, $V \equiv U[U^T U]^{-1} = UQ^{-1}$ は U の一般逆行列で $V^T U = I$ を満たす. また, $\lambda_{ij}, \theta_{ij}$ はスカラー, $\Lambda_{ij}, \tilde{\Lambda}_{ij}, \Theta_{ij}, \tilde{\Theta}_{ij}$ は $K \times K$ 行列で, オーダーパラメータより求めることができる (付録 2 参照). また, Λ_{ij} と $\tilde{\Lambda}_{ij}$ は互いに表裏の関係になって, 片方が求まれば λ_{ij}, Q を用いて残りも求まる. Θ_{ij} と $\tilde{\Theta}_{ij}$ についても同様である. 二通りの表現を用いるのは, それぞれに計算式が簡潔になる利点があるためである.

J の更新則は,

$$\Delta J_i = -\frac{\eta}{\sigma^2 N} \sum_{k=1}^K G_{ik}^{-1} [\delta_k \xi] \quad (35)$$

なので, G^{-1} を含む他のパラメータが固定されているとして, $N \rightarrow \infty$ の極限を計算すると

$$\begin{aligned} \frac{\partial R}{\partial \alpha} &= -\eta \left[\theta \langle \delta y^T \rangle_{\{z\}} + \Xi R \right] \\ \frac{\partial S}{\partial \alpha} &\equiv -\eta \left[\theta \langle \delta x^T \rangle_{\{z\}} + \Xi Q \right] \\ \frac{\partial Q}{\partial \alpha} &= \frac{\partial S}{\partial \alpha} + \frac{\partial S^T}{\partial \alpha} + \eta^2 \theta \langle \delta \delta^T \rangle_{\{z\}} \theta^T \end{aligned} \quad (36)$$

となる. 但し,

$$\Xi \equiv [\Xi_{ij}]_{i,j=1,\dots,K}, \quad \Xi_{ij} \equiv \sum_{k,l=1}^K \langle \delta_k x_l \rangle_{\{z\}} [\Theta_{ik}^T]_{lj} \quad (37)$$

とした.

3.4 適応自然勾配法

ここでは, 次の二つの近似を使う. 1) J と \hat{G}^{-1} の更新に用いる ξ は同一であるが, そうすると二回目以降の更新で J と \hat{G}^{-1} に共通の偏りが生じて, 微分を求める際にパラメータを固定できなくなるため, 別々の独立な ξ を用いるものとして計算する. 本来は J 及び \hat{G}^{-1} のダイナミクスを同時に議論すべきであるが, この仮定により, それぞれの時間微分が別個の平均操作で求まるようになる. 尚, η, ρ が共に小さい時はこの近似が成立する.

2) 二点目は, 式 (12) の平均操作は分母に確率変数があり難しいため, 代わりに \hat{G} のダイナミクスを求めてその逆行列で代用する点である. \hat{G} の現在時刻 α より微小時間 τ 後の値は,

$$\begin{aligned} \hat{G}(\alpha + \tau) &= \left\{ 1 - \frac{\rho}{N} \right\}^{\tau N} \hat{G}(\alpha) \\ &+ \sum_{k=0}^{\tau N - 1} \frac{\rho}{N} \left\{ 1 - \frac{\rho}{N} \right\}^k \frac{1}{\sigma^2} [\nabla f][\nabla f]^T \end{aligned} \quad (38)$$

で, J が固定されているものとして $N \rightarrow \infty$ により,

$$\hat{G}(\alpha + \tau) = e^{-\rho\tau} \hat{G}(\alpha) + \{1 - e^{-\rho\tau}\} G \quad (39)$$

$$\frac{\partial \hat{G}}{\partial \alpha} = \lim_{\tau \rightarrow 0} \frac{1 - e^{-\rho\tau}}{\tau} [-\hat{G} + G] = \rho [-\hat{G} + G] \quad (40)$$

となる. \hat{G}^{-1} を $[\hat{G}]^{-1}$ で代用するので,

$$\begin{aligned} \frac{\partial \hat{G}^{-1}}{\partial \alpha} &= -\hat{G}^{-1} \frac{\partial \hat{G}}{\partial \alpha} \hat{G}^{-1} \\ &= \rho \left[\hat{G}^{-1} - \hat{G}^{-1} G \hat{G}^{-1} \right] \end{aligned} \quad (41)$$

とできる. ρ が小さい時はこの近似が成立する.

次に, $\hat{G}^{-1} = [\hat{G}_{ij}^{-1}]_{i,j=1,\dots,K}$ をオーダーパラメータで表すことを考える. 今, \hat{G}^{-1} が先の $\theta, \Theta, \tilde{\Theta}$ に相当する $\omega, \Omega, \tilde{\Omega}$ を用いて常に表せると仮定する. 即ち,

$$\begin{aligned} \omega &\equiv [\omega_{ij}]_{i,j=1,\dots,K} \\ \Omega &\equiv [\Omega_{ij}]_{i,j=1,\dots,K}, \quad \tilde{\Omega} \equiv [\tilde{\Omega}_{ij}]_{i,j=1,\dots,K} \end{aligned} \quad (42)$$

$$\Omega_{ij} = [\tilde{\Omega}_{ij} - \omega_{ij} I] Q^{-1}, \quad \tilde{\Omega}_{ij} = \Omega_{ij} Q + \omega_{ij} I \quad (43)$$

を用いて,

$$\hat{G}_{ij}^{-1} = \sigma^2 [\omega_{ij} I + U \Omega_{ij} U^T] \quad (44)$$

$$= \sigma^2 \left[\omega_{ij} [I - UV^T] + U \tilde{\Omega}_{ij} V^T \right] \quad (45)$$

とできると仮定する. すると, $\frac{\partial \hat{G}^{-1}}{\partial \alpha}$ も $\omega, \Omega, \tilde{\Omega}$ を用いて表せる.

$$\begin{aligned} \frac{\partial \hat{G}_{ij}^{-1}}{\partial \alpha} &= \rho \sigma^2 \left[[\omega - \omega \lambda \omega]_{ij} [I - UV^T] \right. \\ &\left. + U [\tilde{\Omega} - \tilde{\Omega} \Lambda \tilde{\Omega}]_{[ij]} V^T \right] \end{aligned} \quad (46)$$

但し, $[\bullet]_{[ij]}$ は, i, j 成分ではなく, i, j ブロック ($K \times K$ 行列) を指すものとする. もし, J のダイナミクスが固定 ($\frac{\partial J}{\partial \alpha} = 0$) で, \hat{G}^{-1} のみを更新していくなら, $\omega, \Omega, \tilde{\Omega}$ のダイナミクスは

$$\frac{\partial \omega_{ij}}{\partial \alpha} = \rho [\omega - \omega \lambda \omega]_{ij} \quad (47)$$

$$\frac{\partial \tilde{\Omega}_{ij}}{\partial \alpha} = \rho [\tilde{\Omega} - \tilde{\Omega} \Lambda \tilde{\Omega}]_{[ij]} \quad (48)$$

となるが, 実際には学習則に従って U, V も更新されていくので, $\frac{\partial \tilde{\Omega}}{\partial \alpha}$ はそれを考慮しなければならない. 式 (46) の右辺と式 (45) の右辺を微分した式が等しい

$$\begin{aligned} \rho \sigma^2 \left[[\omega - \omega \lambda \omega]_{ij} [I - UV^T] + U [\tilde{\Omega} - \tilde{\Omega} \Lambda \tilde{\Omega}]_{[ij]} V^T \right] \\ = \sigma^2 \left[\frac{\partial \omega_{ij}}{\partial \alpha} [I - UV^T] - \omega_{ij} \left[\frac{\partial U}{\partial \alpha} V^T + U \frac{\partial V^T}{\partial \alpha} \right] \right. \\ \left. + \frac{\partial U}{\partial \alpha} \tilde{\Omega}_{ij} V^T + U \tilde{\Omega}_{ij} \frac{\partial V^T}{\partial \alpha} + U \frac{\partial \tilde{\Omega}_{ij}}{\partial \alpha} V^T \right] \end{aligned} \quad (49)$$

ことから, 両辺に左右から V^T, U を掛けることで

$$\begin{aligned} \frac{\partial \tilde{\Omega}_{ij}}{\partial \alpha} &= \rho [\tilde{\Omega} - \tilde{\Omega} \Lambda \tilde{\Omega}]_{[ij]} \\ &+ \Omega_{ij} \left[\frac{\partial Q}{\partial \alpha} - \frac{\partial S}{\partial \alpha} \right] - Q^{-1} \frac{\partial S^T}{\partial \alpha} \Omega_{ij} Q \end{aligned} \quad (50)$$

を得る. 式 (47) と (50) を元の式 (49) に代入して, 右辺 - 左辺を計算すると,

$$[I - H] \frac{\partial U}{\partial \alpha} \Omega_{ij} Q V^T + U \Omega_{ij} Q \frac{\partial V^T}{\partial \alpha} [I - H] \quad (51)$$

となり, これは何れの項も 0 になるので, 確かに成立していることが分かる. 但し, $H \equiv UV^T$ は, N 次元ベクトルを U の列ベクトルで張る空間に垂直に落とす直交射影行列, $[I - H]$ は, その垂直成分のみを抽出する行列である.

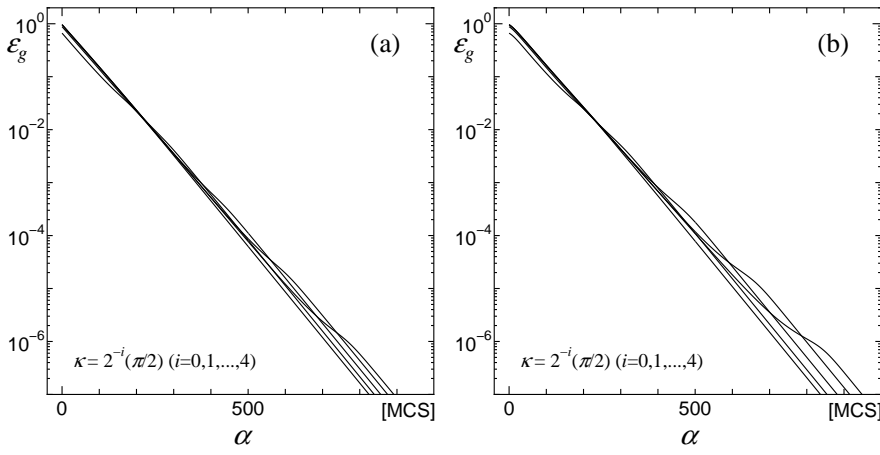


図2 理論計算による汎化誤差の時間変化．(a) 自然勾配法，
(b) 適応自然勾配法 (近似 1(別個の ξ), $2(\widehat{G}^{-1}$ を $[\widehat{G}]^{-1}$ で代用) を使用)．

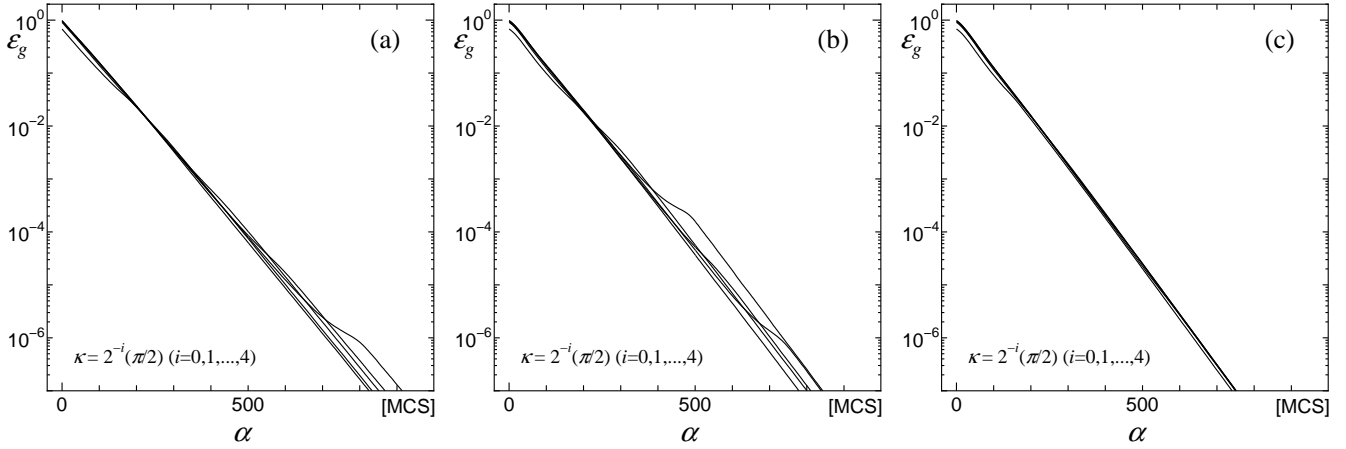


図3 $N = 500$ での汎化誤差の時間変化．(a) 自然勾配法，
(b) 適応自然勾配法 (近似 1(別個の ξ) 使用)，(c) 適応自然勾配法 (近似なし)．

$\omega, \Omega, \widetilde{\Omega}$ の初期値は， $\omega = I, \Omega = 0$ とすれば，先の初期値 $\widehat{G}^{-1} = \sigma^2 I$ を満たすことができる．従って， \widehat{G}^{-1} が常に $\omega, \Omega, \widetilde{\Omega}$ を用いて表せることが分かる．

$\frac{\partial R}{\partial \alpha}, \frac{\partial Q}{\partial \alpha}$ は，前述の自然勾配法のダイナミクスで， $\theta, \Theta, \widetilde{\Theta}$ を $\omega, \Omega, \widetilde{\Omega}$ に置き換えることで同様に計算できる．

4. 結果

自然勾配法と適応自然勾配法の理論を，時間幅 $\alpha = 0.1$ で Runge-Kutta 法で数値計算した結果を示す．中間層ユニット数は生徒，教師共に 2 個 ($K = M = 2$) とし， J の更新割合を小さく $\eta = 10^{-2}$ ， \widehat{G}^{-1} の更新割合を η よりも大きくして $\rho = 10^{-1}$ とした． $\sigma^2 = 1$ は汎化誤差以外には影響しない．生徒，教師の各重みベクトル J_i, B_i については長さを 1 とし，教師間角度 $\kappa \equiv \arccos(\frac{T_{1,2}}{\sqrt{T_{1,1}T_{2,2}}})$ は，幾つか変化させたものを試した．

$$\kappa = 2^{-i}\{\pi/2\} \quad (i = 0, 1, \dots, 4) \quad (52)$$

生徒同士の相関及び対応する教師との相関は， $Q_{1,2} = R_{1,1} = R_{2,2} = \frac{1}{\sqrt{500}}$ ，相手の教師に対する相関は $R_{1,2} = R_{2,1} = \frac{1}{\sqrt{500}} - \frac{2}{\sqrt{500}} \sin \frac{\kappa}{2}$ とした．これらは， J の初期値を $\mathcal{N}(0, \frac{1}{N}I)$ でランダムに設定すると，

$$\langle Q_{ii} \rangle_{\{J\}} = 1$$

$$\sqrt{\langle \{Q_{ij}\}^2 \rangle_{\{J\}}} = \frac{1}{\sqrt{N}} \quad (i \neq j)$$

$$\sqrt{\langle \{R_{ij}\}^2 \rangle_{\{J\}}} = \frac{1}{\sqrt{N}} \quad (53)$$

$$\sqrt{\langle \{R_{ik} - R_{jk}\}^2 \rangle_{\{J\}}} = \frac{2}{\sqrt{N}} \sin \frac{\kappa}{2} \quad (i \neq j)$$

となることから， $N = 500$ を想定した値を用いた．

図 2 に，汎化誤差 ϵ_g の時間変化を示す．(a) の自然勾配法に比べて (b) 適応自然勾配法でもプラトーは殆ど生じず，教師の相関 κ にもあまり影響されていないことが分かる．総じて， $\frac{\eta}{\rho} \ll 1$ の時はプラトーが生じず，そうでない時はプラトーがこの値に応じて増大する傾向があった．これは， $\frac{\eta}{\rho}$ が小さければ， J の変化 (即ち G^{-1} の変化) に \widehat{G}^{-1} の更新が十分ついて行けることによると考えられる．

次に，オリジナルの漸化式 (12) を用いた， $N = 500$ の実シミュレーションを行い，理論を検証した．初期条件は， $J \sim \mathcal{N}(0, \frac{1}{N}I)$ とし，他は理論計算と同様とした．結果を図 3 に示す．(a) は自然勾配法，(b)，(c) は適応自然勾配法で，(b) は ξ を J と \widehat{G}^{-1} 用に別個に用意したもの，(c) は本来の近似を用いない手法である．(b)，(c) 共に図 2(b) の理論値と大きくは異なっていないことが分かる．図 3(c) と図 3(b) の相違は，近似 1) の ξ を別個に用意したことによる

もの、図3(b)と図2(b)の相違は、近似2)の \widehat{G}^{-1} を $[\widehat{G}]^{-1}$ で代用したことによるものと考えられる。

ρ の値が大きい時、適応自然勾配法の理論値は自然勾配法とほぼ同じであったが、近似1,2による二つの相違は共に拡大する傾向があり、近似を用いないシミュレーションは、理論値の倍以上の学習スピードを示した。

5. 結 論

ソフトコミティーマシンのオンライン学習について、統計力学的手法を用いて適応自然勾配法の性能を評価した。その結果、 $\frac{\rho}{\rho}$ を十分小さく設定すれば、自然勾配法とほぼ同等の性能を持つことが分かった。また、理論値を求める際に二つの近似を用い、シミュレーションで近似による誤差の程度を評価した。

付 録

1. ガウス積分

$z \sim \mathcal{N}(\mathbf{0}, \mathbf{C}), g(x) \equiv \text{erf}(x/\sqrt{2}), g'(x) = \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}}$ とすると、以下の積分は解析的に解くことができる [1]。

$$\begin{aligned} \langle g(z_i)g(z_j) \rangle_{\{z\}} &= \frac{2}{\pi} \arcsin \frac{C_{ij}}{\sqrt{\{C_{ii}+1\}\{C_{jj}+1\}}} \\ \langle g'(z_i)g'(z_j)z_k \rangle_{\{z\}} &= \frac{2}{\pi} \frac{c_1}{\{C_{ii}+1\}\sqrt{|M_{ij}|}} \\ \langle g'(z_i)g'(z_j)g'(z_k)g'(z_l) \rangle_{\{z\}} &= \frac{4}{\pi^2} \frac{1}{\sqrt{|M_{ij}|}} \arcsin \frac{c_4}{\sqrt{c_2 c_3}} \end{aligned} \quad (\text{A}\cdot 1)$$

但し、

$$\begin{aligned} M_{ij} &\equiv \begin{bmatrix} C_{ii}+1 & C_{ij} \\ C_{ij} & C_{jj}+1 \end{bmatrix} \\ c_1 &\equiv \begin{vmatrix} C_{ii}+1 & C_{ij} \\ C_{ik} & C_{jk} \end{vmatrix}, c_2 \equiv \begin{vmatrix} M_{ij} & C_{ik} \\ C_{ik} & C_{jk} & C_{kk}+1 \end{vmatrix} \\ c_3 &\equiv \begin{vmatrix} M_{ij} & C_{il} \\ C_{il} & C_{jl} & C_{ll}+1 \end{vmatrix}, c_4 \equiv \begin{vmatrix} M_{ij} & C_{il} \\ C_{ik} & C_{jk} & C_{kl} \end{vmatrix} \end{aligned} \quad (\text{A}\cdot 2)$$

これらを用いて、

$$\begin{aligned} \epsilon_g &= \frac{1}{2\sigma^2} \sum_{i,j=1}^{K+M} a_i a_j \langle g(z_i)g(z_j) \rangle_{\{z\}} \\ \langle \delta_i z_j \rangle_{\{z\}} &= \sum_{k=1}^{K+M} a_k \langle g'(z_i)g'(z_k)z_j \rangle_{\{z\}} \\ \langle \delta_i \delta_j \rangle_{\{z\}} &= \sum_{k,l=1}^{K+M} a_k a_l \langle g'(z_i)g'(z_j)g'(z_k)g'(z_l) \rangle_{\{z\}} \end{aligned} \quad (\text{A}\cdot 3)$$

が求まる。

2. 自然勾配法

自然勾配法の $\lambda, \Lambda, \tilde{\Lambda}, \theta, \Theta, \tilde{\Theta}$ は、上記の M_{ij} を用いて以下のように求めることができる [4]。

$$\lambda_{ij} = \frac{2}{\pi} \frac{1}{\sqrt{|M_{ij}|}} \quad (\text{A}\cdot 4)$$

$$\Lambda_{ij} = -\lambda_{ij} [\mathbf{e}_i \ \mathbf{e}_j] \mathbf{M}_{ij}^{-1} [\mathbf{e}_i \ \mathbf{e}_j]^T \quad (\text{A}\cdot 5)$$

$$\theta = \lambda^{-1}, \quad \tilde{\Theta} = \tilde{\Lambda}^{-1} \quad (\text{A}\cdot 6)$$

但し、 $\mathbf{e}_i \in \mathbb{R}^K$ は第*i*成分のみ1の単位ベクトルである。 $\tilde{\Lambda}$ 及び、 $\tilde{\Theta}$ は、式(28,32)を用いて求まる。

文 献

- [1] D. Saad and A. Solla: Phys. Rev. E **52** (1995) 4225.
- [2] K. Fukumizu and S. Amari: Neural Networks **13** (2000) 317.
- [3] S. Amari: Neural Comput. **10** (1998) 251.
- [4] H. H. Yang and S. Amari: Neural Comput. **10** (1998) 2137.
- [5] M. Rattray and D. Saad: Phys. Rev. E **59** (1999) 4523.
- [6] M. Inoue, H. Park and M. Okada: J. Phys. Soc. Jpn. **72** (2003) (*in press*).
- [7] S. Amari, H. Park and K. Fukumizu: Neural Comput. **12** (2000) 1399.
- [8] H. Park, S. Amari and K. Fukumizu: Neural Networks **13** (2000) 755.