

## アンサンブル学習の解析と拡張

内田 真人<sup>†</sup>      塩谷 浩之<sup>††</sup>

伊達 惇<sup>††</sup>

Analysis and Extension of Ensemble Learning

Masato UCHIDA<sup>†</sup>, Hiroyuki SHIOYA<sup>††</sup>,  
and Tsutomu DA-TE<sup>††</sup>,

<sup>†</sup> NTT サービスインテグレーション基盤研究所, 武蔵野市

NTT Service Integration Laboratories, Musashino-shi, 180-8585 Japan

<sup>††</sup> 北海道大学大学院工学研究科, 札幌市

Graduate School of Engineering, Hokkaido University,  
Sapporo-shi, 060-8628 Japan

あらまし 汎化能力向上の目的から、個々に学習された複数の学習機械の出力を適当に重み付けしたものを予測値として用いるというアンサンブル学習が提案されている。本論文ではこの学習法の特性を、負の対数ゆ度関数を損失関数とする枠組みを用いて解析する。また、この枠組みを拡張し、アンサンブル学習を意味づける。

キーワード アンサンブル学習, 平均予測誤差,  $\alpha$ -ダイバージェンス

### 1. まえがき

個々に学習された複数の学習機械が与えられたとき、汎化能力の高い予測機械をどのように得るかという問題は重要である。アンサンブル学習は、与えられた機械の出力を適当に重み付けして予測値とすることでこの問題への対応を試みる手法であり、次のように定式化されている [4]。

$M$  個の異なる学習機械を考え、各々、入力  $x = (x_1, \dots, x_m)^T \in \mathbb{R}^m$  に対して  $f(x; \theta_i) \in \mathbb{R}$  を出力するものとする ( $i = 1, \dots, M$ )。ただし、 $\theta_i = (\theta_i^{(1)}, \dots, \theta_i^{(k_i)})^T \in \mathbb{R}^{k_i}$  は学習機械  $f(x; \theta_i)$  の修正可能なパラメータであり、 $f(x; \theta_i)$  は  $\theta_i$  に関して高階微分可能とする。また、入力  $x$  と望ましい出力  $y$  からなるサンプルの組  $(x, y)$  が、ある確率分布  $p_*(x, y) (= q(x)p_*(y|x))$  に従って互いに独立に  $\sum_{i=1}^M n_i$  個観測されたとし、これを  $D_{n_i} = \{(x_1^{(i)}, y_1^{(i)}), \dots, (x_{n_i}^{(i)}, y_{n_i}^{(i)})\}$  と書く。このとき、2乗誤差関数を損失関数とするアンサンブル学習とは

$$\hat{\theta}_i = \arg \min_{\theta_i} \sum_{(x, y) \in D_{n_i}} (y - f(x; \theta_i))^2 \quad (1)$$

により与えられる

$$\bar{f}(x; \hat{\theta}, \beta) = \sum_{i=1}^M \beta_i f(x; \hat{\theta}_i) \quad (2)$$

を予測機械として用いることをいう。ただし

$$\begin{aligned} \sum_{i=1}^M \beta_i &= 1, \quad \beta_i > 0 \\ \beta &= (\beta_1, \dots, \beta_M)^T \in \mathbb{R}^M \\ \theta &= (\theta_1^T, \dots, \theta_M^T)^T \in \mathbb{R}^{\sum_{i=1}^M k_i} \end{aligned}$$

である。

本論文では上記の定式化を特別な場合として含む枠組みを用い、アンサンブル学習をパラメータ推定の立場から漸近論的に解析する。また、この枠組みに基づいてアンサンブル学習のアルゴリズム構造を検討する。

### 2. 準備

#### 2.1 アンサンブル学習の別表現

ある関数  $g(x) \in \mathbb{R}$  を用いて、 $x$  に対する  $y$  の条件付き確率分布

$$p_G(y|g(x)) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2}(y - g(x))^2\right\}$$

を定めると ( $\sigma$  は正定数)、式 (1) は

$$\hat{\theta}_i = \arg \min_{\theta_i} \left\{ - \sum_{(x, y) \in D_{n_i}} \log p_G(y|f(x; \theta_i)) \right\} \quad (3)$$

と書き換えられ、式 (2) は

$$p_G(y|\bar{f}(x; \hat{\theta}, \beta)) = \frac{\prod_{i=1}^M p_G(y|f(x; \hat{\theta}_i))^{\beta_i}}{\int_{\mathbb{R}} \prod_{i=1}^M p_G(y|f(x; \hat{\theta}_i))^{\beta_i} dy} \quad (4)$$

と恒等的に等しい [1]。

本論文では上記の対応関係に注目し、負の対数ゆ度関数を損失関数とするアンサンブル学習の定式化と解析を行う。これは、2乗誤差関数を損失関数とするアンサンブル学習を特別な場合として含む。次節では、そのための補題を用意する。

#### 2.2 補題の準備

集合  $\mathcal{Z}$  上の確率分布全体を  $\mathcal{P}(\mathcal{Z})$  とおく。すなわち

$$\mathcal{P}(\mathcal{Z}) \stackrel{\text{def}}{=} \left\{ p|p: \mathcal{Z} \rightarrow \mathbb{R}, p(z) \geq 0 (\forall z \in \mathcal{Z}), \right.$$

$$\left. \int_{\mathcal{Z}} p(z) dz = 1 \right\}$$

とおく．ここで， $p_i(z) (\in \mathcal{P}_i(\mathcal{Z}) \subset \mathcal{P}(\mathcal{Z}))$  を用いて ( $i = 1, \dots, M$ )，新たな確率分布  $\bar{p}_\beta(z) (\in \mathcal{P}(\mathcal{Z}))$  を

$$\bar{p}_\beta(z) \stackrel{\text{def}}{=} \frac{\prod_{i=1}^M p_i(z)^{\beta_i}}{\int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz}$$

と定義する．ただし

$$\int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz < \infty$$

$$\sum_{i=1}^M \beta_i = 1, \quad \beta_i > 0$$

$$\beta = (\beta_1, \dots, \beta_M)^T \in \mathbb{R}^M$$

とする．このとき，以下の補題が成立する．ただし， $D(\cdot \parallel \cdot)$  は Kullback ダイバージェンスを表す．

[補題 1] 上の定義のもとで

$$D(p \parallel \bar{p}_\beta) = \sum_{i=1}^M \beta_i D(p \parallel p_i) + \log \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz$$

が成り立つ． □

補題 1 から，仮に，添字  $i$  によらず  $D(p \parallel p_i)$  の値が同じであれば， $D(p \parallel \bar{p}_\beta)$  の  $\beta$  に関する最小化は

$$-\min_{\beta} \log \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz$$

を求めることに相当することがわかる．この量は Chernoff 情報量と呼ばれる [1]．

[補題 2] 上の定義のもとで

$$-\log \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz = \sum_{i=1}^M \beta_i D(\bar{p}_\beta \parallel p_i) \geq 0$$

が成り立つ． □

また，次が知られている [1]．

[補題 3] 上の定義のもとで

$$\bar{p}_\beta(z) = \arg \min_{p(z) \in \mathcal{P}(\mathcal{Z})} \sum_{i=1}^M \beta_i D(p \parallel p_i)$$

が成り立つ．別表現として

$$\sum_{i=1}^M \beta_i D(\bar{p}_\beta \parallel p_i) \leq \sum_{i=1}^M \beta_i D(\tilde{p} \parallel p_i)$$

が成り立つ．ここで  $\tilde{p}(z)$  は  $\mathcal{P}(\mathcal{Z})$  の任意の要素を表す． □

### 3. 解 析

#### 3.1 問題設定

2.2 の仮定を満たす確率分布  $p_i(z) \stackrel{\text{def}}{=} p(z; \theta_i)$  ( $\in \mathcal{P}_i(\mathcal{Z}) \subset \mathcal{P}(\mathcal{Z})$ ) について ( $i = 1, \dots, M$ )

$$\begin{aligned} \bar{p}_\beta(z) &\stackrel{\text{def}}{=} \bar{p}(z; \theta, \beta) \\ &\stackrel{\text{def}}{=} \frac{\prod_{i=1}^M p(z; \theta_i)^{\beta_i}}{\int_{\mathcal{Z}} \prod_{i=1}^M p(z; \theta_i)^{\beta_i} dz} \end{aligned} \quad (5)$$

とおく．ただし

$$\theta_i = (\theta_i^{(1)}, \dots, \theta_i^{(k_i)})^T \in \mathbb{R}^{k_i}$$

$$\theta = (\theta_1^T, \dots, \theta_M^T)^T \in \mathbb{R}^{\sum_{i=1}^M k_i}$$

であり， $p(z; \theta_i)$  は  $\theta_i$  に関して高階微分可能とする．また，ある確率分布  $p_*(z)$  に従って互いに独立に  $\sum_{i=1}^M n_i$  個のサンプルが観測されたとし，これを  $D_{n_i} = \{z_1^{(i)}, \dots, z_{n_i}^{(i)}\}$  と書く．このとき，2.1 に対応させて，負の対数ゆう度関数を損失関数とするアンサンブル学習を

$$\hat{\theta}_i = \arg \min_{\theta_i} \left\{ - \sum_{z \in D_{n_i}} \log p(z; \theta_i) \right\} \quad (6)$$

によって与えられる  $\bar{p}(z; \hat{\theta}, \beta)$  を求めることであると定義し，その汎化能力を平均予測誤差

$$E_{\hat{\theta}}[D(p_*(z) \parallel \bar{p}(z; \hat{\theta}, \beta))] \quad (7)$$

で評価する．ただし， $E$  は添字の従う確率分布に関する期待値を表す．また

$$\mathcal{P}_i(\mathcal{Z}) = \{p(z; \theta_i) \mid \theta_i \in \mathbb{R}^{k_i}\}$$

$$\mathcal{P}_1(\mathcal{Z}) \subset \dots \subset \mathcal{P}_M(\mathcal{Z})$$

$$\exists \theta_i^* \text{ s.t. } p_*(z) = p(z; \theta_i^*) \in \mathcal{P}_i(\mathcal{Z})$$

とする．更に

$$\check{\Theta}_i = \{\theta \mid \text{正射影} : \theta \mapsto \theta_i, \theta_i \in \mathbb{R}^{k_i}, \theta \in \mathbb{R}^{k_M}\}$$

とし， $\theta_i \in \mathbb{R}^{k_i}$  と  $\check{\theta}_i \in \check{\Theta}_i$  を同一視する．ただし， $\check{\theta}_i$  のうち  $\theta_i$  と対応しない部分は  $\theta_j^*$  と一致するものと仮定する ( $\forall j > i$ )．以下では，記法の簡単のため  $\check{\theta}_i$  と  $\theta_i$  を区別せずに用いる．例えば， $\theta_i + \theta_j$  は  $\check{\theta}_i + \check{\theta}_j$  を意味する ( $i < j$ )．また，この仮定のもとで  $\theta_i^* = \theta^*$  となる ( $i = 1, \dots, M$ )．

### 3.2 平均予測誤差の上界と下界

式 (7) を直接解析することは困難であるため、その上界と下界を評価する。

[補題 4] 式 (6) で表される  $\hat{\theta}_i$  について、 $n_i$  が十分大きいとき ( $i = 1, \dots, M$ )、適当な正則条件のもとで

$$\begin{aligned} E_{\hat{\theta}}[D(p(z; \hat{\theta}) \| p(z; \hat{\theta}_i))] \\ = \frac{k_i}{2n_i} - \frac{k_i}{\sum_{l=1}^M n_l} + \frac{1}{2} \frac{\sum_{j=1}^M n_j k_j}{(\sum_{l=1}^M n_l)^2} \end{aligned}$$

が成り立つ。ただし

$$\hat{\theta} \stackrel{\text{def}}{=} \frac{\sum_{i=1}^M n_i \hat{\theta}_i}{\sum_{i=1}^M n_i}$$

である。

補題 1, 2, 3, 4 より、次が成り立つ。

[定理 1]  $n_i$  が十分大きいとき ( $i = 1, \dots, M$ )、適当な正則条件のもとで

$$\begin{aligned} \frac{\sum_{i=1}^M \beta_i k_i}{\sum_{j=1}^M n_j} - \frac{1}{2} \frac{\sum_{i=1}^M n_i k_i}{(\sum_{j=1}^M n_j)^2} \\ \leq E_{\hat{\theta}}[D(p_*(z) \| \bar{p}(z; \hat{\theta}, \beta))] \\ \leq \sum_{i=1}^M \frac{\beta_i k_i}{2n_i} \end{aligned}$$

が成り立つ。

[系 1]  $n = n_i$ ,  $k = k_i$  のとき ( $i = 1, \dots, M$ )、 $n$  が十分大きければ、適当な正則条件のもとで

$$\frac{k}{2Mn} \leq E_{\hat{\theta}}[D(p_*(z) \| \bar{p}(z; \hat{\theta}, \beta))] \leq \frac{k}{2n}$$

が成り立つ。

### 3.3 考 察

系 1 の下界は  $Mn$  個のサンプルを用いたときの平均予測誤差、上界は  $n$  個のサンプルを用いたときの平均予測誤差を表している [2]。したがって、 $Mn$  個のサンプルが一括で与えられている場合には、アンサンブル学習は有効ではなく、 $n$  個のサンプルを用いて学習した予測機械のみが  $M$  個与えられている場合にはアンサンブル学習は有効であるということがわかる。この様子は [4] において実験的に示されていた。また、一般にアンサンブル学習の有効性の根拠は、補題 2 で与えられる量が正值となることにあることが容易にわかる。一方、新たなサンプル  $D_{n'} = \{z'_1, \dots, z'_{n'}\}$  が与えられた場合には

$$\hat{\beta} = \arg \min_{\beta} \left\{ - \sum_{z \in D_{n'}} \log \bar{p}(z; \hat{\theta}, \beta) \right\} \quad (8)$$

によって  $\beta$  を設定することができるが、本論文では、これに関する解析結果は与えられていない。

### 4. $\alpha$ -ダイバージェンスへの拡張

前節までに扱ってきたアンサンブル学習のアルゴリズムの本質は

$$\hat{p}_i(z) \stackrel{\text{def}}{=} \arg \min_{p(z) \in \mathcal{P}(Z)} D(p_* \| p) \quad (9)$$

$$\bar{p}_{\beta}(z) \stackrel{\text{def}}{=} \arg \min_{p(z) \in \mathcal{P}(Z)} \sum_{i=1}^M \beta_i D(p \| \hat{p}_i) \quad (10)$$

$$\hat{\beta} \stackrel{\text{def}}{=} \arg \min_{\beta} D(p_* \| \bar{p}_{\beta}) \quad (11)$$

と表される 3 段階の最小化操作に帰着される。なぜならば、式 (1), (3), (6) は式 (9) に、式 (2), (4), (5) は式 (10) に、更に、式 (8) は式 (11) に対応するからである。

ところで、最近では [5] に見られるように、学習問題における  $\alpha$ -ダイバージェンスの重要性が論じられている。そこで本節では式 (9), (10), (11) の最小化操作が、 $\alpha$ -ダイバージェンスを用いた場合に自然に拡張できることを明らかにする。ただし、 $\alpha$ -ダイバージェンスは次のように定式化されている ( $0 < \alpha < 1$ ) [3]。

$$D^{(\alpha)}(p \| q) = \frac{1}{\alpha(1-\alpha)} \left\{ 1 - \int_{\mathcal{Z}} p(z)^{\alpha} q(z)^{1-\alpha} dz \right\}$$

まず、補題 3 は次のように拡張される。

[定理 2] 上の定義のもとで

$$\begin{aligned} \bar{p}_{\beta}^{(\alpha)}(z) &\stackrel{\text{def}}{=} \frac{\{\sum_{i=1}^M \beta_i p_i(z)^{1-\alpha}\}^{\frac{1}{1-\alpha}}}{\int_{\mathcal{Z}} \{\sum_{i=1}^M \beta_i p_i(z)^{1-\alpha}\}^{\frac{1}{1-\alpha}} dz} \\ &= \arg \min_{p(z) \in \mathcal{P}(Z)} \sum_{i=1}^M \beta_i D^{(\alpha)}(p \| p_i) \end{aligned}$$

が成り立つ。

[系 2] 上の定義のもとで

$$\begin{aligned} \lim_{\alpha \rightarrow 1} \bar{p}_{\beta}^{(\alpha)}(z) &= \bar{p}_{\beta}(z) \\ \lim_{\alpha \rightarrow 0} \bar{p}_{\beta}^{(\alpha)}(z) &= \sum_{i=1}^M \beta_i p_i(z) \end{aligned}$$

が成り立つ。

また、補題 1, 2 は次のように拡張される。

[定理 3] 上の定義のもとで

$$\begin{aligned} C^{(\alpha)} D^{(\alpha)}(p \| \bar{p}_{\beta}^{(\alpha)}) \\ = \sum_{i=1}^M \beta_i D^{(\alpha)}(p \| p_i) - \sum_{i=1}^M \beta_i D^{(\alpha)}(\bar{p}_{\beta}^{(\alpha)} \| p_i) \end{aligned}$$

が成り立つ．ただし

$$C(\alpha) \stackrel{\text{def}}{=} \sum_{i=1}^M \beta_i \int_{\mathcal{Z}} \{\bar{p}_\beta^{(\alpha)}(z)\}^\alpha \{p_i(z)\}^{1-\alpha} dz$$

である．また

$$\lim_{\alpha \rightarrow 1} C(\alpha) = \lim_{\alpha \rightarrow 0} C(\alpha) = 1$$

が成り立つ．  $\square$

以上から， $\alpha$ -ダイバージェンスを用いたアンサンブル学習のアルゴリズムは次のように定義するのが自然である．

$$\hat{p}_i^{(\alpha)}(z) \stackrel{\text{def}}{=} \arg \min_{p(z) \in \mathcal{P}_i(\mathcal{Z})} D^{(\alpha)}(p_* \| p)$$

$$\bar{p}_\beta^{(\alpha)}(z) \stackrel{\text{def}}{=} \arg \min_{p(z) \in \mathcal{P}(\mathcal{Z})} \sum_{i=1}^M \beta_i D^{(\alpha)}(p \| \hat{p}_i^{(\alpha)})$$

$$\hat{\beta}^{(\alpha)} \stackrel{\text{def}}{=} \arg \min_{\beta} D^{(\alpha)}(p_* \| \bar{p}_\beta^{(\alpha)})$$

## 5. むすび

本論文では，負の対数ゆう度関数を損失関数とするアンサンブル学習を定式化し，その平均予測誤差の上界と下界を求めた．この結果，すべてのサンプルが一括で与えられている場合には，アンサンブル学習は有効ではなく，学習済みの学習機械のみが与えられている場合にはアンサンブル学習は有効であるということがわかった．また，一般にアンサンブル学習の有効性の根拠は，補題 2 で与えられる量が正值となることにあるがわかった．これらの結果は 2 乗誤差関数を損失関数とする場合にも成立する．また，アンサンブル学習のアルゴリズムの本質が，式 (9)，(10)，(11) で表される，Kullback ダイバージェンスに関する 3 段階の最小化操作にあることから，これを  $\alpha$ -ダイバージェンスを用いる場合に拡張した．今後の課題としては， $\beta$  の推定についての考察や，これまでに筆者らが行った研究 [6] との関連性についての考察などが挙げられる．

## 文 献

- [1] T.M. Cover and J.A. Thomas, Elements of Information Theory, Wiley-Interscience publication, America, 1991.
- [2] 坂本慶行, 石黒真木夫, 北川源四郎, 情報量統計学, 情報科学講座 A・5・4, 共立出版, 1983.
- [3] 甘利俊一, 長岡浩司, 情報幾何の方法, 岩波講座応用数学 6 [対象 12], 岩波書店, 1993.
- [4] 上田修功, 中野良平, “アンサンブル学習における汎化誤差解析” 信学論 (D-II), vol.J80-D-II, no.9, pp.2512-2521, Sept. 1997.
- [5] 松山泰男, “ $\alpha$ -EM アルゴリズムとその基本的性質” 信学論

- (D-I), vol.J82-D-I, no.12, pp.1347-1358, Dec. 1999.  
 [6] 内田真人, 塩谷浩之, 伊達 惇, “非ベイズの付加項を用いた多層パーセプトロンの学習” 信学論 (D-II), vol.J83-D-II, no.6, pp.1572-1576, June 2000.

## 付 録

### 1. 補題 1 の証明

定義より

$$\begin{aligned} D(p \| \bar{p}) &= \int_{\mathcal{Z}} p(z) \log \prod_{i=1}^M \left( \frac{p(z)}{p_i(z)} \right)^{\beta_i} dz \\ &+ \int_{\mathcal{Z}} p(z) \left\{ \log \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z')^{\beta_i} dz' \right\} dz \\ &= \sum_{i=1}^M \beta_i D(p \| p_i) + \log \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz \end{aligned}$$

が成立する．

### 2. 補題 2 の証明

定義より

$$\begin{aligned} -\log \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz &= -\int_{\mathcal{Z}} \bar{p}_\beta(z) \left\{ \log \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z')^{\beta_i} dz' \right\} dz \\ &= \int_{\mathcal{Z}} \bar{p}_\beta(z) \log \prod_{i=1}^M \left( \frac{\bar{p}_\beta(z)}{p_i(z)} \right)^{\beta_i} dz \\ &= \sum_{i=1}^M \beta_i D(\bar{p}_\beta \| p_i) \end{aligned}$$

が成立する．不等号は， $\beta$  の定義と Kullback ダイバージェンスの非負性より明らか．

### 3. 補題 4 の証明

$-\log p(z; \hat{\theta}_i)$  を  $\hat{\theta}$  の周りで Taylor 展開し，高次の項を無視することで

$$\begin{aligned} -\log p(z; \hat{\theta}_i) &= -\log p(z; \tilde{\theta}) - (\hat{\theta}_i - \tilde{\theta})^T \nabla \log p(z; \tilde{\theta}) \\ &- \frac{1}{2} (\hat{\theta}_i - \tilde{\theta})^T \nabla \nabla^T \log p(z; \tilde{\theta}) (\hat{\theta}_i - \tilde{\theta}) \end{aligned}$$

が得られる．ただし

$$\nabla = \left( \frac{\partial}{\partial \theta^{(1)}}, \dots, \frac{\partial}{\partial \theta^{(k_M)}} \right)^T$$

であり,  $\nabla \nabla^T \log p(\mathbf{z}; \boldsymbol{\theta})$  は  $\log p(\mathbf{z}; \boldsymbol{\theta})$  の Hessian 行列 ( $\partial^2 \log p(\mathbf{z}; \boldsymbol{\theta}) / \partial \theta^{(i)} \partial \theta^{(j)}$ ) を表す. よって

$$\begin{aligned} & D(p(\mathbf{z}; \tilde{\boldsymbol{\theta}}) \| p(\mathbf{z}; \hat{\boldsymbol{\theta}}_i)) \\ &= \frac{1}{2} (\hat{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}})^T G(\tilde{\boldsymbol{\theta}}) (\hat{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}) \\ &= \frac{1}{2} \left\{ 1 - \frac{n_i}{\sum_{l=1}^M n_l} \right\}^2 (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}^*)^T G(\tilde{\boldsymbol{\theta}}) (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}^*) \\ &\quad + \frac{1}{2} \sum_{j=1, j \neq i}^M \left\{ \frac{n_j}{\sum_{l=1}^M n_l} \right\}^2 \\ &\quad \times (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}^*)^T G(\tilde{\boldsymbol{\theta}}) (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}^*) \\ &\quad + \sum_{i \neq j} N_{ij} (\hat{\boldsymbol{\theta}}_i - \boldsymbol{\theta}^*)^T G(\tilde{\boldsymbol{\theta}}) (\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}^*) \quad (\text{A.1}) \end{aligned}$$

が成り立つ. ただし,  $G(\boldsymbol{\theta})$  は Fisher 情報行列  $E_{\mathbf{Z}}[\nabla \log p(\mathbf{z}; \boldsymbol{\theta}) \nabla \log p(\mathbf{z}; \boldsymbol{\theta})^T]$  であり,  $N_{ij}$  は  $\hat{\boldsymbol{\theta}}_i, \hat{\boldsymbol{\theta}}_j$  に依存しない定数である. また, 各  $n_i$  ( $i = 1, \dots, M$ ) が十分大きいとき, 適当な正則条件のもとで

$$\hat{\boldsymbol{\theta}}_i \sim N\left(\boldsymbol{\theta}^*, \frac{1}{n_i} G(\boldsymbol{\theta}^*)^{-1}\right)$$

であり [2], かつ, 各  $\hat{\boldsymbol{\theta}}_i$  ( $i = 1, \dots, M$ ) は独立であることから

$$\tilde{\boldsymbol{\theta}} \sim N\left(\boldsymbol{\theta}^*, \frac{1}{\sum_{i=1}^M n_i} G(\boldsymbol{\theta}^*)^{-1}\right)$$

となる. そこで, 式 (A.1) において高次の項を無視して  $G(\tilde{\boldsymbol{\theta}}) = G(\boldsymbol{\theta}^*)$  と置き換える. 以上から

$$\begin{aligned} & E_{\tilde{\boldsymbol{\theta}}} [D(p(\mathbf{z}; \tilde{\boldsymbol{\theta}}) \| p(\mathbf{z}; \hat{\boldsymbol{\theta}}_i))] \\ &= E_{\hat{\boldsymbol{\theta}}_1} \cdots E_{\hat{\boldsymbol{\theta}}_M} [D(p(\mathbf{z}; \tilde{\boldsymbol{\theta}}) \| p(\mathbf{z}; \hat{\boldsymbol{\theta}}_i))] \\ &= \frac{1}{2} \left\{ 1 - \frac{n_i}{\sum_{l=1}^M n_l} \right\}^2 \frac{k_i}{n_i} \\ &\quad + \frac{1}{2} \sum_{j=1, j \neq i}^M \left\{ \frac{n_j}{\sum_{l=1}^M n_l} \right\}^2 \frac{k_j}{n_j} \\ &= \frac{k_i}{2n_i} - \frac{k_i}{\sum_{l=1}^M n_l} + \frac{1}{2} \frac{\sum_{j=1}^M n_j k_j}{(\sum_{l=1}^M n_l)^2} \end{aligned}$$

が得られる.

#### 4. 定理 1 の証明

$n_i$  が十分大きいとき, 適当な正則条件のもとで

$$E_{\hat{\boldsymbol{\theta}}_i} [D(p_*(\mathbf{z}) \| p(\mathbf{z}; \hat{\boldsymbol{\theta}}_i))] = \frac{k_i}{2n_i}$$

である ( $i = 1, \dots, M$ ) [2]. これと, 補題 2 より上界が示される. 更に, 補題 1, 3, 4 を用いることで下界が示される.

#### 5. 定理 2 の証明

Lagrange 乗数  $\lambda$  を用いて汎関数  $L(q)$  を

$$L(q) \stackrel{\text{def}}{=} \sum_{i=1}^M \beta_i D^{(\alpha)}(q \| p_i) + \lambda \int_{\mathcal{Z}} q(\mathbf{z}) d\mathbf{z}$$

と定義する. ただし,  $q(\mathbf{z}) \in \mathcal{P}(\mathcal{Z})$  である. このとき

$$\begin{aligned} \frac{\partial L(q)}{\partial q(\mathbf{z})} &= -\frac{q(\mathbf{z})^{\alpha-1}}{1-\alpha} \left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z})^{(1-\alpha)} \right\} + \lambda \\ &= 0 \end{aligned}$$

を満たす  $q(\mathbf{z})$  が求めるものである. ここで

$$\int_{\mathcal{Z}} q(\mathbf{z}) d\mathbf{z} = 1$$

となるように  $\lambda$  を選ぶことで定理が得られる.

#### 6. 定理 3 の証明

前半を証明する. まず

$$\begin{aligned} & \left[ \int_{\mathcal{Z}} \left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z})^{1-\alpha} \right\}^{\frac{1}{1-\alpha}} d\mathbf{z} \right]^{1-\alpha} \\ &= \int_{\mathcal{Z}} \frac{\left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z})^{1-\alpha} \right\}^{\frac{1}{1-\alpha}}}{\left[ \int_{\mathcal{Z}} \left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z}')^{1-\alpha} \right\}^{\frac{1}{1-\alpha}} d\mathbf{z}' \right]^\alpha} \\ &\quad \times \left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z})^{1-\alpha} \right\}^{\frac{\alpha-1}{1-\alpha}} \left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z})^{1-\alpha} \right\} d\mathbf{z} \\ &= \int_{\mathcal{Z}} \{ \bar{p}_\beta^{(\alpha)}(\mathbf{z}) \}^\alpha \left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z})^{1-\alpha} \right\} d\mathbf{z} \\ &= \sum_{i=1}^M \beta_i \int_{\mathcal{Z}} \{ \bar{p}_\beta^{(\alpha)}(\mathbf{z}) \}^\alpha \{ p_i(\mathbf{z}) \}^{1-\alpha} d\mathbf{z} \\ &\stackrel{\text{def}}{=} C(\alpha) \quad (\text{A.2}) \end{aligned}$$

が成り立つ. ここで, 定義と式 (A.2) を用いることで

$$\begin{aligned} & C(\alpha) D^{(\alpha)}(p \| \bar{p}_\beta^{(\alpha)}) \\ &= \frac{1}{\alpha(1-\alpha)} \left\{ 1 - \int_{\mathcal{Z}} p(\mathbf{z})^\alpha \right\} \end{aligned}$$

$$\times \left\{ \sum_{i=1}^M \beta_i p_i(\mathbf{z})^{1-\alpha} \right\} d\mathbf{z} \Bigg\} \\ - \frac{1}{\alpha(1-\alpha)} \{1 - C(\alpha)\}$$

$$= \sum_{i=1}^M \beta_i D^{(\alpha)}(p \| p_i) - \sum_{i=1}^M \beta_i D^{(\alpha)}(\bar{p}_\beta^{(\alpha)} \| p_i)$$

が示される．後半は明らか．

(平成 12 年 10 月 17 日受付 , 13 年 1 月 9 日再受付)