

アンサンブル学習モデルにおける重み付けに関する考察

内田 真人[†](正員) 塩谷 浩之^{††}(正員)

A Study on Assignment of Weight Parameters in Ensemble Learning Model

Masato UCHIDA[†] and Hiroyuki SHIOYA^{††}, *Regular Members*

[†] 日本電信電話株式会社 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, NTT Corporation, 3-9-11 Midori-cho, Musashino-shi, 180-8585 Japan

^{††} 室蘭工業大学工学部情報工学科, 室蘭市

Department of Computer Science and Systems Engineering, Muroran Institute of Technology, 27-1 Mizumoto-cho, Muroran-shi, 050-8585 Japan

あらまし アンサンブル学習とは、個々に学習された複数の学習機械の予測結果の重み付け平均を最終的な予測結果として用いる手法である。この学習法は、Kullback divergence に関する 3 段階の最小化操作として理解できる場合がある。本論文では、この枠組みにおける重み付けに関する考察を行う。

キーワード アンサンブル学習, 重み付け, 初期値変動法, 三角不等式

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。アンサンブル学習とは、与えられた機械の予測結果の重み付け平均を最終的な予測結果とすることでこの問題への対応を試みる手法である。この種の学習法の手続きは様々に設定可能であると思われるが、学習機械の混合に関する確率モデルを適当に設定することで、学習法の手続きを明確に表現できる場合がある [4]。この概要を以下に示す。

集合 \mathcal{Z} ($\subset \mathbb{R}^d$) 上の確率分布全体を

$$\mathcal{P}(\mathcal{Z}) \stackrel{\text{def}}{=} \left\{ p \mid p: \mathcal{Z} \rightarrow \mathbb{R}, p(z) > 0 (\forall z \in \mathcal{Z}), \int_{\mathcal{Z}} p(z) dz = 1 \right\}$$

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

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

と定義する。ただし

$$\int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\beta_i} dz < \infty \quad (2)$$

$$\sum_{i=1}^M \beta_i = 1, \quad \beta_i > 0 \quad (3)$$

$$\beta = (\beta_1, \dots, \beta_M)^T \in \mathbb{R}^M \quad (4)$$

とする。

このとき、以下のような Kullback divergence に関する 3 段階の最小化操作を考える。これは、2 乗誤差関数を損失関数とし、与えられた機械の出力の重み付け平均を予測値とするという学習法 [3] の自然な一般化である [4]。ただし、 $D(\cdot \parallel \cdot)$ は Kullback divergence を表し、 p_* は推定の対象となる真の分布である。

$$\hat{p}_i(z) \stackrel{\text{def}}{=} \arg \min_{p(z) \in \mathcal{P}_i(\mathcal{Z})} D(p_* \parallel p) \quad (5)$$

$$\bar{\hat{p}}_\beta(z) \stackrel{\text{def}}{=} \arg \min_{p(z) \in \mathcal{P}(\mathcal{Z})} \sum_{i=1}^M \beta_i D(p \parallel \hat{p}_i) \quad (6)$$

$$\hat{\beta} \stackrel{\text{def}}{=} \arg \min_{\beta \in \mathbb{R}^M} D(p_* \parallel \bar{\hat{p}}_\beta) \quad (7)$$

上記の手続きにより求められた $\bar{\hat{p}}_\beta$ が、最終的に推定された確率分布であり、各 \hat{p}_i ($i = 1, \dots, M$) が事前に推定された確率分布である。この枠組みを用いることにより、与えられたサンプルを分割し個々の学習機械を生成するという手法に基づくアンサンブル学習は、平均予測誤差の意味で有効ではないという事実が導かれる ([4] の系 1)。

一方で [4] においては \bar{p}_β のパラメータ β に関する議論が十分になされていない。例えば、 \bar{p}_β が確率分布となるためには式 (2) の条件が与えられていれば十分であるにもかかわらず、この条件に加えて、パラメータ β が重みとしての意味をもつための条件である式 (3) が追加されている点に関する考察が十分ではない。更に、式 (7) により定まる $\hat{\beta}$ がこの条件を満たすかどうかとも全く自明ではない。

次章では、与えられたサンプルをすべて用い個々の学習機械を生成する、という手法に基づくアンサンブル学習の性質を重み β の最適化の観点から述べ、更に、式 (3) がアンサンブル学習の結果とどのようにかわってくるのかを調べる。

2. 主結果

式 (7) により定まる $\hat{\beta}$ の特徴は、式 (1) により定められる確率モデルが指数型分布族に属することから

$$D(p\|\bar{p}_\beta) = D(p\|p_i) - D(\bar{p}_\beta\|p_i), \quad (i = 1, \dots, M) \quad (8)$$

のように与えられることが知られている [1].

式 (8) は, アンサンブル学習において, 与えられたすべてのサンプルを用いた上で, パラメータをこう配法により決定する際の初期値を変えて得られた複数の学習機械を利用すること (初期値変動法 [3]) が, [4] の系 1 によって得られた, アンサンブル学習に関する否定的な結果を解決する一つの方法であることを示している. なぜならば, $Q(\mathcal{Z}) = \{p_1(z), \dots, p_M(z)\}$ を, 異なる初期値を用いて学習した M 個の学習機械とし, 一般性を失うことなく

$$p_M(z) = \arg \min_{q(z) \in Q(\mathcal{Z})} D(p\|q) \quad (9)$$

としたとき, 最適な重み $\hat{\beta}$ が得られたとすれば, 式 (8) より, $D(p\|\bar{p}_\beta)$ は $D(p\|p_M)$ よりも $D(\bar{p}_\beta\|p_M)$ だけ小さくなるからである. この様子は [3] において実験的に示されていた.

[注意] サンプルを分割するという手法においては, β の最適化を行っても, 上記のような性質は得られない ([4] の系 1 参照). □

ところで, このときの $\hat{\beta}$ は式 (3) の条件を必ずしも満たさない. つまり, $\hat{\beta}$ を用いた場合には, 重み付き平均をとるという意味を与えられない場合があるということである. しかし, 次が成り立つ.

[定理] 式 (8) における β が一意に定まるとき

$$0 \leq \hat{\beta}_i \leq 1, \quad (i = 1, \dots, M) \quad (10)$$

が成り立つための必要十分条件は

$$D(p\|p_M) + D(p_M\|p_i) \geq D(p\|p_i)$$

である ($i = 1, \dots, M$). ただし, 一般性を失うことなく

$$D(p\|p_M) = \min_{i=1, \dots, M} D(p\|p_i)$$

と仮定した. □

(証明については付録を参照)

上記の定理から, 式 (10) を満たすためには, 各 i について, p_i 及び p_M の p (真の分布) に対する近似精度の差 $D(p\|p_i) - D(p\|p_M)$ が $D(p_M\|p_i)$ より小さいことが要請される. これは, p_i の近似精度が悪い場合には, 式 (10) が成り立たないことを意味する.

なお, 式 (10) が成り立たない場合, ある j, k について $\beta_k > 1, \beta_j < 0$ が成り立つ. これは, $\beta_j (< 0)$ により p_j の高密度部を低密度部に変換し, 更に, $\beta_k (> 0)$ により p_k の高密度部を際立たせることになる. 無論, 上記の p_k は複数の場合もある.

3. むすび

本論文では, アンサンブル学習における重み β の最適化に関する考察を加えた. その結果, 最適な重み $\hat{\beta}$ が得られたとすれば, [4] の系 1 により結論されるアンサンブル学習に関する否定的な結果を解決する一つの方法として, 初期値変動法は有効であるということがわかった. この結果や [4] の系 1 の主張は, 本論文で用いたアンサンブル学習の理論モデルの特徴に基づき導出されたものである. したがって, 式 (5), (6), (7) により与えられる本論文の枠組みは, [3] の数値実験により示されている結果を説明するために適した理論モデルの一つであるといえる.

一方, $\hat{\beta}$ については, それがアンサンブル学習における重みとしての意味をもつための必要十分条件を Kullback divergence に関する三角不等式により与えた.

今後の課題としては, 近似精度の低い学習機械を用いた場合のアンサンブル学習における汎化誤差解析などが挙げられる.

文 献

- [1] 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.
- [2] 大関信雄, 青柳雅計, 不等式, 槇書店, 1976.
- [3] 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.
- [4] 内田真人, 塩谷浩之, 伊達 惇 “アンサンブル学習の解析と拡張” 信学論 (D-II), vol. J84-D-II, no.7, pp.1537-1542, July 2001.

付 録

定理の証明

一般性を失うことなく

$$\beta_M = 1 - \sum_{i=1}^{M-1} \beta_i$$

としてよい. このとき

$$\begin{aligned} \frac{\partial}{\partial \beta_j} D(p\|\bar{p}_\beta) &= D(p\|p_j) - D(p\|p_M) \\ &+ \int_{\mathcal{Z}} \bar{p}_\beta(z) \log \frac{p_j(z)}{p_M(z)} dz \quad (A.1) \end{aligned}$$

が成り立つ . よって

$$\begin{aligned} & \frac{\partial^2}{\partial \beta_j^2} D(p \parallel \bar{p}_\beta) \\ &= \int_{\mathcal{Z}} \bar{p}_\beta(z) \left\{ \log \frac{p_j(z)}{p_M(z)} \right\}^2 dz \\ & \quad - \left\{ \int_{\mathcal{Z}} \bar{p}_\beta(z) \log \frac{p_j(z)}{p_M(z)} dz \right\}^2 \\ &= \int_{\mathcal{Z}} \bar{p}_\beta(z) \left\{ \log \frac{p_j(z)}{p_M(z)} \right. \\ & \quad \left. - \int_{\mathcal{Z}} \bar{p}_\beta(z') \log \frac{p_j(z')}{p_M(z')} dz' \right\}^2 dz \\ & \geq 0 \end{aligned}$$

となるので , $D(p \parallel \bar{p}_\beta)$ は β_j に関して下に凸である ($j = 1, \dots, M-1$) . また , 式 (A.1) より

$$\begin{aligned} \frac{\partial}{\partial \beta_j} D(p \parallel \bar{p}_\beta) \Big|_{\beta_j=1} &= D(p \parallel p_j) - D(p \parallel p_M) \\ & \quad + D(p_j \parallel p_M) \geq 0 \end{aligned}$$

が成り立つ ($j = 1, \dots, M-1$) . したがって

$$0 \leq \hat{\beta}_j \leq 1$$

であるための必要十分条件は

$$\begin{aligned} \frac{\partial}{\partial \beta_j} D(p \parallel \bar{p}_\beta) \Big|_{\beta_j=0} &= D(p \parallel p_j) - D(p \parallel p_M) \\ & \quad - D(p_M \parallel p_j) \leq 0 \end{aligned}$$

である ($j = 1, \dots, M-1$) . この不等式は $j = M$ のときには自明に成り立つ . 一方

$$\begin{aligned} 0 \leq \hat{\beta}_j \leq 1, \quad (j = 1, \dots, M-1) \\ \rightarrow \hat{\beta}_M = 1 - \sum_{i=1}^{M-1} \hat{\beta}_i \leq 1 \end{aligned}$$

が成り立つ . よって

$$\begin{aligned} 0 \leq \hat{\beta}_j \leq 1, \quad (j = 1, \dots, M-1) \\ \rightarrow \hat{\beta}_M = 1 - \sum_{i=1}^{M-1} \hat{\beta}_i \geq 0 \end{aligned}$$

が成り立てば証明は完了する . これを背理法で示すために $\hat{\beta}_M < 0$ と仮定する . このとき , 一般性を失うことなく $\beta_j \neq 0$ ($j = 1, \dots, M-1$) とすると

$$\begin{aligned} & \int_{\mathcal{Z}} \prod_{i=1}^M p_i(z)^{\hat{\beta}_i} dz \\ & \geq (1 - \hat{\beta}_M) \int_{\mathcal{Z}} \left\{ \prod_{i=1}^{M-1} p_i(z)^{\hat{\beta}_i} \right\}^{\frac{1}{1-\hat{\beta}_M}} dz \\ & \quad + \hat{\beta}_M \int_{\mathcal{Z}} p_M(z) dz \\ & \geq \int_{\mathcal{Z}} \left\{ \prod_{i=1}^{M-1} p_i(z)^{\hat{\beta}_i} \right\}^{\frac{1}{1-\hat{\beta}_M}} dz \\ & \quad + \hat{\beta}_M \int_{\mathcal{Z}} \left\{ p_M(z) - \sum_{i=1}^{M-1} \frac{\hat{\beta}_i}{1-\hat{\beta}_M} p_i(z) \right\} dz \\ & = \int_{\mathcal{Z}} \prod_{i=1}^{M-1} p_i(z)^{\frac{\hat{\beta}_i}{1-\hat{\beta}_M}} dz \end{aligned}$$

が成り立つ . ここで , $a > 0, b > 0, \frac{1}{u} + \frac{1}{v} = 1$, ($u \neq 0$) のとき

$$\begin{aligned} u < 1 \quad \text{ならば} \quad a^{\frac{1}{u}} b^{\frac{1}{v}} &\geq \frac{a}{u} + \frac{b}{v} \\ u > 1 \quad \text{ならば} \quad a^{\frac{1}{u}} b^{\frac{1}{v}} &\leq \frac{a}{u} + \frac{b}{v} \end{aligned}$$

であり , 等号は $a = b$ のときかつそのときに限る [2] ことを用いた . また

$$\begin{aligned} & \sum_{i=1}^M \hat{\beta}_i D(p \parallel p_i) \\ &= \sum_{i=1}^{M-1} \frac{\hat{\beta}_i}{1-\hat{\beta}_M} D(p \parallel p_i) + \hat{\beta}_M D(p \parallel p_M) \\ & \quad + \sum_{i=1}^{M-1} \left(\hat{\beta}_i - \frac{\hat{\beta}_i}{1-\hat{\beta}_M} \right) D(p \parallel p_i) \\ & \geq \sum_{i=1}^{M-1} \frac{\hat{\beta}_i}{1-\hat{\beta}_M} D(p \parallel p_i) + \hat{\beta}_M D(p \parallel p_M) \\ & \quad + \sum_{i=1}^{M-1} \left(\hat{\beta}_i - \frac{\hat{\beta}_i}{1-\hat{\beta}_M} \right) D(p \parallel p_M) \\ &= \sum_{i=1}^{M-1} \frac{\hat{\beta}_i}{1-\hat{\beta}_M} D(p \parallel p_i) \end{aligned}$$

が成り立つ . よって

$$\begin{aligned} \tilde{\beta}_j &= \frac{\hat{\beta}_j}{1-\hat{\beta}_M}, \quad j = 1, \dots, M-1 \\ \tilde{\beta} &= (\tilde{\beta}_1, \dots, \tilde{\beta}_{M-1})^T \end{aligned}$$

とおけば [4] の補題 1 と $\tilde{\beta}$ の一意性より

$$D(p\|\bar{p}_{\hat{\beta}}) < D(p\|\bar{p}_{\beta})$$

$\hat{\beta}_M \geq 0$ である .

が成り立つ . これは , $\hat{\beta}$ の定義に矛盾する . すなわち ,

(平成 14 年 7 月 10 日受付 , 12 月 19 日再受付)
