

モデル平均理論の新展開

劉慶豊*

小樽商科大学

平成21年7月31日

概要

本論文はモデル選択とその発展系であるモデル平均(model averaging)に関する既存の研究結果を概観する。モデル平均の意味合いとモデル平均の応用及び未解決問題に関して説明する。さらに、モデル平均に関連する未解決問題に関する研究の将来的な方向を述べる。論文の中で実際のデータを利用した予測の例を示す。

1 はじめに

確率モデルをもとにデータ解析を行う際に、候補となるモデル群から一つ最良と思われるモデルを選ぶ手段として、モデル選択理論は盛んに研究されてきた。実証研究の主流として、AIC(Akaike [1973]) やBIC(Schwarz [1978])などの情報量基準を利用した選択法が広く利用されている。一方で、モデル選択に伴う不確実性を考慮して、ベイズ的な立場をとる研究者が異なったフィロソフィーに基づいて、ベイジアン・モデル平均の手法と理論を構築し発展させてきた。また、ここ数年、非ベイズ的な研究者は非ベイジアン・モデル平均理論を構築し始めた。本稿はモデル選択とモデル平均の理論を概観し、異なるフィロソフィーに基づく異なる手法の意味合いを明らかにする。さらに、この分野の未解決問題と研究における将来の発展方向の可能性について論じる。最後に実証例を挙げて、実証研究に適用するときの注意点を述べる。

第2章で、代表的な情報量基準であるAICやBICなどを紹介する。そして、第3章でモデル平均の意味合いを述べ、幾つかの例を挙げて、モデル平均の

* 小樽商科大学商学部, Email:qliu@res.otaru-uc.ac.jp

推定量の性質に関して述べる。第4章で未解決問題に関して論じ、第5章でモデル平均を利用した幾つかの研究例を紹介した上で、実際のデータを用いて、ファイナンスへの応用を示す。第6章は結論である。

1.1 情報量基準によるモデル選択

統計学者や計量経済学者はモデルを構成して、データの特性をできるだけ正確に説明しようとする。一般的にはデータを生成する真の確率分布を近似するための候補として複数のモデルが考えられる。それらのモデルを評価するためにさまざまな情報量基準が提案されている。なかでも、AICやBICなどは最も広く利用されている。

1.2 赤池情報量基準AIC

AICの基本となっているものはカルバックーライブラー距離である。確率変数 y の標本 $\mathbf{y} = (y_1, y_2, \dots, y_n)$ が観測されたとする。 y の真の密度関数を $g(y)$ としてパラメトリックなモデル $f(y, \boldsymbol{\theta})$ で近似しようとする。ただし、 $\boldsymbol{\theta} \in \Theta$ で $\Theta \subseteq \mathbb{R}^K$ である。 $f(y, \boldsymbol{\theta})$ による近似の良さを計るために一つの基準として、Kullback and Leibler [1951]がカルバックーライブラー距離を提案した。その定義は以下のようになる、

$$\begin{aligned} KL(g, f) &= E_g \left(\log \frac{g(y)}{f(y, \boldsymbol{\theta})} \right) \\ &= \int_{-\infty}^{\infty} \{ \log g(y) - \log f(y, \boldsymbol{\theta}) \} g(y) dy. \end{aligned} \quad (1)$$

$\log(\cdot)$ は自然対数を表す。カルバックーライブラー距離はモデルで真の分布を近似するときの一種のリスク関数として考えられる。

$KL(g, f)$ の中の $g(y)$ はモデルに依存しないため、モデルの良さを評価する場合、 $\log g(y)$ の部分を無視できる。原理的には平均対数尤度と呼ばれる

$$E_g (\log f(y, \boldsymbol{\theta})) = \int_{-\infty}^{\infty} \log f(y, \boldsymbol{\theta}) g(y) dy \quad (2)$$

を最大にするようなモデルが最適なモデルとなる。ここで、大数の法則を利用すれば

$$\lim \frac{1}{n} \sum_{i=1}^n \log f(y_i, \boldsymbol{\theta}) = E_g (\log f(y, \boldsymbol{\theta})) \quad (3)$$

が分かる。そのため、対数尤度関数を

$$l_n(\boldsymbol{\theta}) = \sum_{i=1}^n \log f(y_i, \boldsymbol{\theta}) \quad (4)$$

と定義し、 $\hat{\boldsymbol{\theta}}$ を最尤推定量として、 $\frac{1}{n}l_n(\hat{\boldsymbol{\theta}})$ が平均対数尤度 $E_g(\log f(y, \hat{\boldsymbol{\theta}}))$ の一つの推定量として考えられるが、この推定量はバイアスを持っていることが知られている（証明は小西・北川 [2004]を参照されたい）。そこで、バイアスを修正して、モデル選択の基準としてAICが構築された。AICは

$$\text{AIC} = -2l_n(\hat{\boldsymbol{\theta}}) + 2p \quad (5)$$

と定義される。ただし、 p はパラメーター $\boldsymbol{\theta}$ の次元で、平均対数尤度を推定するときのバイアスの補正として働く。 $2p$ はパラメーターの数、すなわちモデルの複雑さに対するペナルティとして解釈されることもある。モデルを選択する際、AICを最小にするモデルを選ぶ。

1.3 他の情報量基準

AICと並んでベイジアン情報量基準(BIC)も広く利用されている。他に、AICの修正版として竹内 [1976]とStone [1977]によって考案された情報量基準TICもよく知られている。

BICの背景となっているのはベイジアン統計学である。事後確率 $P(M|\mathbf{y})$ が最大になるようなモデル M を選ぶというのがその発想である。導出の詳細はClaeskens and Hjort [2008]を参照されたい。事後確率の近似的な評価値としてBICが定義される。

$$\text{BIC} = -2l_n(\hat{\boldsymbol{\theta}}) + (\log n)p \quad (6)$$

AICと似た形になっているが、そのペナルティ $(\log n)p$ は n が大きくなるときAICのペナルティより大きくなる。 n が大きくなるに連れAICはより複雑なモデルを選びがちであるが、BICは真のモデルよりパラメーター数の少ないモデルを選ぶ傾向がある。そして、真のモデルが候補となるモデル族の中に入つていれば、サンプル数 n が無限大に近づくにつれ、BICが真のモデルを選ぶ確率は1に収束する。すなわち、BICはモデル選択の情報量基準として一致性を持つ。その一方では、AICは一致性を持たない。

小西・北川 [2004]によれば、AICを計算する際、モデルがデータを生成する確率分布と一致するときに限って、 p が平均対数尤度を推定するときのバイ

アスの補正項になることが保障される。そうではないとき, p が正しいバイアスの補正項とならない。竹内 [1976]は正確な補正項を導出して, 想定したモデル $f(y, \boldsymbol{\theta})$ の中に真の確率分布を含まれていない場合で利用できる情報量基準(TIC) を構築した。

$$TIC = -2l_n(\hat{\boldsymbol{\theta}}) + 2Tr\left(\hat{I}\hat{J}^{-1}\right) \quad (7)$$

\hat{I} と \hat{J} はそれぞれ

$$I_0 = E_g\left[\frac{\partial \log f(y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \frac{\partial \log f(y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}'}\right] \quad (8)$$

$$J_0 = -E_g\left[\frac{\partial^2 \log f(y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right] \quad (9)$$

の推定量で, 一致推定量 $\hat{\boldsymbol{\theta}}$ を I_0 と J_0 のサンプル版に代入することで得られる。モデルはデータを生成する確率分布と一致しないとき, AICが厳密的には正確な補正項を使っていないが, 想定するモデルが真の確率分布に近い場合, 近似的にはAICをモデル選択の手法として利用することが可能で, 実際にも広く使われている。

また, 線形回帰モデルの情報量基準としてMallows [1973]はMallows' C_p

$$C_p = RSS + 2p\hat{\sigma}^2 - n\hat{\sigma}^2 \quad (10)$$

を定義した。 C_p はリスク関数Mean Squared Error(MSE)の推定量となっている。それを最小にするように, モデルに入れる説明変数を選ぶ。 C_p は線形モデルにしか適用できないが, 非線形モデルは線形化してからの適用が可能である。

2 モデル平均に関して

実証研究で, AICやBICなどの情報量基準を利用してモデル選択を行う際, 選ばれたモデルをそのまま真のモデルとして推論が展開されることが多い。しかし, AICやBICなどの情報量基準によるモデル選択は一定の不確実性をもたらす。その不確実性が推定量の性質に影響を及ぼすことを考慮しないと間違った推論に繋がる。そして, その推論のもとで計算した信頼区間や行った検定の結果が誤ったものになる。この様な問題が選択後推定量 (Post Model

Selection Estimator, PMSE) の問題と呼ばれる。モデル平均はこのような PMSE の問題を考慮に入れた手法である。

モデル平均はモデル選択を拡張して得られた手法である。候補となるモデル族を $\mathcal{M} = \{M_1, M_2, \dots, M_K\}$ とする。ある母数 μ のモデル平均による推定量の一般表現は

$$\hat{\mu} = \sum_{M_k \in \mathcal{M}} c(M_k) \hat{\mu}_{M_k} \quad (11)$$

で与えられる。 $\hat{\mu}_{M_k}$ はモデル M_k のもとで得られる推定量で、 $c(M_k)$ は各モデルに与えられるウェイトを決める関数である。 $c(M_k)$ は総和が 1 で、非確率的でも確率的でも良い。この表現は情報量基準でモデルを選択して得られる推定量を含んでいる。たとえば、AIC の場合、AIC で選んだモデルを M_{AIC} とする。そのモデルによって得られる推定量は

$$\hat{\mu}_{AIC} = \sum_{M_k \in \mathcal{M}} I(M_k = M_{AIC}) \hat{\mu}_{M_k} \quad (12)$$

となる。一つ注意すべきは、AIC を基準にモデルを選ぶ際、どのモデルが選ばれるかは確率的なイベントであり、 $I(M_k = M_{AIC})$ が確率変数となっていることである。後ほどこの点に関して詳しく論じることにする。

2.1 モデル選択はモデル平均の特殊ケース

AIC で選んだモデルによって得られた推定量の場合、最小の AIC に対応しているモデルにウェイト 1 を与えて、残りのモデルに 0 のウェイトを与えていている。モデル平均による推定量の場合、ウェイトはより一般的で、AIC の場合のウェイトを特殊ケースとして含む。そのため、リスクで評価するとき、リスクに関する最適なウェイトを使えば、必ずモデル平均による推定量のリスクは、AIC の場合のものと等しいかより小さくなる。この意味ではモデル平均の手法を利用することで、単純に AIC などの情報量基準による手法より優れた結果を得る可能性がある。他の特定の一つのモデルしか選ばない伝統的なモデル選択の手法に関しても同様のことと言える。

$\hat{\mu}_{AIC}$ の例や、モデル平均による推定量の一般表現からも分かるように、推定量の性質を導出する際、ウェイトとなる部分の確率的な特性による影響、すなわち、モデル選択に伴う不確実性を考慮しなければならない。例えば、 $\hat{\mu}_{AIC}$ は本質的に M_{AIC} が真のモデルであるという条件のもとで得られた推定量と異なるものである。そのため、AIC でモデル M_{AIC} を選択し、選択されたモデル

を真のモデルとして推定量の性質を導出するのは誤ったやり方である。これはいわゆるモデル選択後推定量 (Post-Model-Selection Estimator, PMSE) の問題である。これから、いくつかのモデル平均の手法を紹介してから、PMSE の問題を論じることにする。

2.2 ベイジアンモデル平均

近年、ベイズ的な枠組みでモデル平均の手法に関する論文が数多く発表されてきた (Draper [1995], Hoeting et al. [1999], Clyde and George [2004])。ベイジアンの枠組みでは、モデルのパラメーターを確率変数として考える。関心のある母数を μ として、各モデルのパラメーター $\boldsymbol{\theta}_k$ の関数とする。その母数の事後分布 $\pi(\mu|\mathbf{y})$ や条件付期待値 $\hat{\mu} = E(\mu|\mathbf{y})$ を推定することが目的である。

これまでの設定と同じく、候補となるモデル族を $\mathcal{M} = \{M_1, M_2, \dots, M_K\}$ とする。ベイズの公式を利用して、モデル平均によって $\pi(\mu|\mathbf{y})$ と $E(\mu|\mathbf{y})$ が導出される。そのためには、二つの事前分布が必要となる。一つは各モデルの事前分布で $P(M_k)$ と記す。もう一つは各モデルのパラメーター $\boldsymbol{\theta}_k$ のそのモデルのもとでのパラメーターの事前分布 $\pi(\boldsymbol{\theta}_k|M_k)$ である。この二つの事前分布を所与として、 \mathbf{y} のモデル M_k のもとでの尤度関数を $L(\mathbf{y}|M_k, \boldsymbol{\theta}_k)$ として、モデル M_k の事後分布は

$$P(M_k|\mathbf{y}) = \frac{P(M_k)\lambda_k}{\sum_{k=1}^K P(M_k)\lambda_k} \quad (13)$$

となる。ただし

$$\lambda_k = \int L(\mathbf{y}|M_k, \boldsymbol{\theta}_k) \pi(\boldsymbol{\theta}_k|M_k) d\boldsymbol{\theta}_k \quad (14)$$

である。さらに μ の事後分布は

$$\pi(\mu|\mathbf{y}) = \sum_{k=1}^K P(M_k|\mathbf{y}) \pi(\mu|M_k, \mathbf{y}) \quad (15)$$

となる。ただし、 $\pi(\mu|M_k, \mathbf{y})$ はモデル M_k が真のモデルである条件のもとでの μ の事後分布であり、事前分布 $\pi(\boldsymbol{\theta}_k|M_k)$ を利用してベイズの公式により解析的に、またはMCMCなどの方法で数値的に推定できる。そして μ の事後期待値は

$$\hat{\mu} = E(\mu|\mathbf{y}) = \sum_{k=1}^K P(M_k|\mathbf{y}) E(\mu|M_k, \mathbf{y}) \quad (16)$$

と導出される。

ベイジアンモデル平均の手法はすでに多くの成果をあげているが、問題点も残されている。たとえば候補となるモデルに与える事前確率を決めるとき、かなり主観的で場当たりな方法が取られる。それにより推定結果が大きく左右される。

2.3 Smoothed-AIC, BIC

Buckland et al. [1997]はモデルのAICおよびBICの値を利用したモデル平均を方法を提案した。この方法はベイジアンモデル平均と密接な関連を持っている。

前述したようにBICは事後確率 $P(M_k|\mathbf{y})$ を評価している。Claeskens and Hjort [2008]によれば近似的に

$$\text{BIC} \approx -2 \log(\lambda_k) \quad (17)$$

が成立する。そこで $P(M_k)$ が k に関して均一であると仮定すれば、(13)式から

$$P(M_k|\mathbf{y}) \approx \frac{\exp(-\text{BIC}_k/2)}{\sum_{k=1}^K \exp(-\text{BIC}_k/2)} \quad (18)$$

であることがわかる。(16)式はモデル平均のウェイトとして

$$c_{BIC}(M_k) = \frac{\exp(-\text{BIC}_k/2)}{\sum_{k=1}^K \exp(-\text{BIC}_k/2)} \quad (19)$$

が一つ妥当な選択肢であることを示唆する。そのウェイトを利用して構築したモデル平均の推定量を

$$\hat{\mu}_{MA-BIC} = \sum_{M_k \in \mathcal{M}} c_{BIC}(M_k) \hat{\mu}_{M_k} \quad (20)$$

と定義できる。この推定量はsmoothed-BIC-based estimatorと呼ばれる。

$\hat{\mu}_{MA-BIC}$ の定義と同じフォームを取って、smoothed-AIC-based estimatorが

$$\hat{\mu}_{MA-AIC} = \sum_{M_k \in \mathcal{M}} c_{AIC}(M_k) \hat{\mu}_{M_k} \quad (21)$$

と定義される。ただし、

$$c_{AIC}(M_k) = \frac{\exp(-\text{AIC}_k/2)}{\sum_{k=1}^K \exp(-\text{AIC}_k/2)} \quad (22)$$

である。モデル族 \mathcal{M} に含まれる各モデルに最尤法を適用すれば、各モデルの推定量とウエイトを計算できるため、この二つのモデル平均の推定量の計算は非常に簡単である。

2.4 Mallows' C_p を用いたモデル平均

Hansen [2007]はMallows' C_p を利用したモデル平均の方法を提案した。モデルは線形モデルとして表現されている。モデル式は

$$y_i = \mu_i + \varepsilon_i \quad (23)$$

$$\mu_i = \sum_{j=1}^{\infty} \theta_j x_{ji} \quad (24)$$

$$E(\varepsilon_i | x_i) = 0 \quad (25)$$

$$E(\varepsilon_i^2 | x_i) = \sigma^2 \quad (26)$$

である。 (y_i, \mathbf{x}_i) , $\mathbf{x}_i = (x_{1i}, x_{2i}, \dots)$, $i = 1, 2, \dots, n$ はランダムなサンプルで、 ε_i が搅乱項である。 $\sum_{j=1}^{\infty} \theta_j x_{ji}$ を非線形またはノンパラメトリックな関数の級数展開として考えることができるので、この方法の適用範囲は線形モデルだけに限定されない。

このモデルのもとで、 \mathbf{x}_i の最初の k 個の要素だけを説明変数とするモデル

$$\mu_i = \sum_{j=1}^k \theta_j x_{ji} \quad (27)$$

を M_k で表す。最小のモデルが M_1 で、最大のモデルを M_K と限定する。モデル平均の推定量は

$$\hat{\boldsymbol{\theta}} = \sum_{k=1}^K \omega_k \begin{pmatrix} \hat{\boldsymbol{\theta}}_k \\ \mathbf{0} \end{pmatrix} \quad (28)$$

と定義される。ただし、 $\hat{\boldsymbol{\theta}}_k = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)'$ は $\boldsymbol{\theta}_k$ のモデル M_k のもとでの推定量で、 $\boldsymbol{\theta}_k = (\theta_1, \theta_2, \dots, \theta_k)'$ である。 ω_k はモデル M_k に与えるウエイトで、 $\omega_k \in [0, 1]$, $\sum_{k=1}^K \omega_k = 1$ を満たす。 K 個のウエイトをまとめて $W = (\omega_1, \omega_2, \dots, \omega_K)'$ と記す。

\mathbf{X} を $\mathbf{x}_i, i = 1, 2, \dots, n$ のデータ行列とする。 $\hat{\mu}(W) = \mathbf{X}\hat{\boldsymbol{\theta}}$ は W をウエイトとして利用して得られた μ のモデル平均の推定量を表す。損失関数を $L_n(W) =$

$(\hat{\mu}(W) - \mu)'(\hat{\mu}(W) - \mu)$ として、Hansen [2007]はMallows' C_p が一定の条件のもとで、損失関数の期待値と定数の和の不偏推定量であることを示した。ただし、Mallows' C_p の定義は行列表現で

$$C_n(W) = \left(Y - X_K \hat{\Theta} \right)' \left(Y - X_K \hat{\Theta} \right) + 2\sigma^2 k^* \quad (29)$$

となる。定義式の中で、 X_K は $k = K$ の場合のモデル、すなわち最大のモデルの説明変数のデータ行列で、 $k^* = \sum_{k=1}^K \omega_k k$ である。さらに、Mallows' C_p 、 $C_n(W)$ を最小にするウエイトを

$$\hat{W}_n = \arg \min_W C_n(W) \quad (30)$$

と定義する。一定の条件のもとで

$$\frac{L_n(\hat{W}_n)}{\inf_W L_n(W)} \xrightarrow{p} 1 \quad (31)$$

が満たされ、 \hat{W}_n は漸近的に損失関数の値を最小にする最適なウエイトとなる。

2.5 PMSE推定量の漸近分布

モデル選択で選んだモデルのもとで得られる推定量、およびモデル平均による推定量は、Post-model-selection(PMSE)推定量と呼ばれている。このように呼ぶのは、それらの推定量はモデル選択やモデル平均のステップによってその漸近的な性質が影響され、推定したモデルが真のモデルであるという仮定のもとでの推定量とは異なることを表現したいためである。モデル選択とモデル平均によってもたらされる不確実性を考慮に入れ、PMSE推定量の漸近分布の導出を行った論文は幾つか存在するが(Hjort and Claeskens [2003], Potscher [2006]とLeeb and Potscher [2008])、ここでは、Hjort and Claeskens [2003]のアプローチに関して簡単に紹介する。

Hjort and Claeskens [2003]はlocal misspecificationのもとで漸近分布の導出を行った。(Potscher [2006]とLeeb and Potscher [2008]も同様に、local misspecificationのもとで漸近分布を導出している。) 真のモデルをサンプル数に依存する密度関数の列で定義している。 $y_i, i = 1, 2, \dots, n$ が確率変数 y のサンプルで、その密度関数は

$$f_{true}(y) = f_n(y) = f(y, \boldsymbol{\theta}_0, \boldsymbol{\gamma}_0 + \boldsymbol{\delta}/\sqrt{n}) \quad (32)$$

とする。 $\boldsymbol{\theta}_0$ と γ_0 は次元が p と K のパラメータベクトルである。最小のモデルは γ に属するパラメーターを含まない、 $f_{narr}(y, \boldsymbol{\theta}) = f(y, \boldsymbol{\theta}, \gamma_0)$ である。最大のモデルはすべてのパラメーターを含むようなモデルで $f_{full}(y, \boldsymbol{\theta}, \gamma)$ と表す。 γ に属するパラメーターの中から幾つかを選んでモデルに入れることによって、最大と最小のモデルの中間にあるサブモデルを構築することができる。 M_{S_j} でサブモデルを表す。 S_j はモデル M_{S_j} に入れた γ の要素の添え字の集合である。 γ の要素のすべての組み合わせを考えると、最大と最小のモデルを含めて、サブモデルは全部で 2^K 種類がある。ある母数 μ を推定したい場合、そのモデル平均推定量は

$$\hat{\mu} = \sum_{j \in 2^K} c(M_{S_j}) \hat{\mu}_{S_j} \quad (33)$$

と定義される。ただし、 $\hat{\mu}_{S_j}$ はモデル M_{S_j} のもとでの推定量である。

Hjort and Claeskens [2003]はこのような設定のもとで、最尤法のテクニックを駆使してモデル平均推定量 $\hat{\mu}$ の漸近分布を導出した。一般的には、 $c(M_{S_j})$ が確率変数となり、導出された分布は標準的な分布に従わず、一種の正規分布の混合分布となる。そして、モデル選択による推定量の漸近分布はその特殊ケースとして含まれる。Hjort and Claeskens [2003]は漸近分布の結果を利用して、モデル選択で得られたモデルを真のモデルとして、信頼空間を導出した場合、信頼空間が過小評価されることを論じた。また、幾つかの異なるモデル平均推定量のリスクの比較を行った。

3 未解決問題

統計的推論に関しては前節で紹介したPMSE推定量の分布を利用して推定量の有限標本の分布を近似し、検定や信頼区間の導出を行うのは一つの自然な考え方である。しかし、このような近似はサンプル数の大きさに関わらず、必要な精度に到達するのは不可能であることがPotscher [2006]やLeeb and Potscher [2008])などにより示されている。

以下ではLeeb and Potscher [2008]の結果を説明する。モデルを行列表現で

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\lambda} + \boldsymbol{\varepsilon} \quad (34)$$

と表す。 \mathbf{X} と \mathbf{Z} はそれぞれ $(n \times k)$ と $(n \times q)$ の非確率的な説明変数の行列である。 $\boldsymbol{\varepsilon} \sim N(0, \sigma^2 \mathbf{I}_n)$ は搅乱項である。パラメーターをまとめて $\boldsymbol{\theta} = (\boldsymbol{\beta}', \boldsymbol{\lambda}')'$ とする。 \mathbf{X} に含まれる説明変数は必ずモデルに入れることにする。モデル選択によ

り, \mathbf{Z} に含まれる説明変数のうち, どれを残すかを決める。そして, 選択されたモデルのもとで, β を含むパラメーターを推定する。 β の推定量 $\hat{\beta}$ の有限標本の分布を $G_{n,\theta,\sigma}(\mathbf{b})$ とする。Leeb and Potscher [2008]によれば, $G_{n,\theta,\sigma}(\mathbf{b})$ の推定量として $\hat{G}_n(\mathbf{b})$ が任意の正の数 δ に関して

$$P_{n,\theta,\sigma}\left(\left|\hat{G}_n(\mathbf{b}) - G_{n,\theta,\sigma}(\mathbf{b})\right| > \delta\right) \xrightarrow{n \rightarrow \infty} 0 \quad (35)$$

を満たすような, すなわち θ に関してポイントワイズに一致性を持つ推定量を構築することができるが, 一方で, ある正の数 η が存在して, その推定量 $\hat{G}_n(\mathbf{b})$ について

$$\sup_{\|\theta\| < \eta} P_{n,\theta,\sigma}\left(\left|\hat{G}_n(\mathbf{b}) - G_{n,\theta,\sigma}(\mathbf{b})\right| > \delta\right) \xrightarrow{n \rightarrow \infty} 1 \quad (36)$$

が成り立つ。これは $\hat{G}_n(\mathbf{b})$ はパラメーター θ に関して一様な一致性を持たないことを意味する。Leeb and Potscher [2008]は, この一様に一致性を持つ推定量が存在しないという結論に不可能性と名付けた。不可能性は以上のモデルの設定のもとだけで成立するのではなく, より一般的なパラメトリックとセミパラメトリックモデルに関する成立する。さらに, 不可能性は殆どのモデル選択の方法, そしてモデル平均に関する成り立つ。

不可能性の結論により, サンプル数を如何に大きくしても, PMSE推定量の有限標本の分布の推定がうまく行かないことが分かる。そのため, 有限標本分布の推定量のもとで推論を行うアプローチが実現できない。そこで, Beran and Dumbgen [1998]とBeran [2000]などはこの問題を解決するための一つの方向を提示している。その大まかな方針は以下のようになっている。まず, 候補となるモデルに対応するリスクを推定する。そして, リスクを最小にするモデルを選ぶ。次は, 選んだモデルのリスクの推定量の漸近分布を導出する。最終的にはリスクの推定量の漸近分布のもとで統計的推論を行う。Beran and Dumbgen [1998]とBeran [2000]はこのアプローチに基づいて推定量を中心とした母数の信頼集合を導出した。

4 モデル平均の応用

モデル選択の手法は各分野で広く利用されている。それとは対照的に, モデル平均の手法の応用はまだ数が少ない。これから, 幾つかのモデル平均を利用した研究を紹介する。その上, 実際のデータを利用して, モデル平均による予測の例を示す。

4.1 既存研究

経済理論を計量経済学の手法を利用して分析するとき、モーメント条件を利用して推定や推論を行うのは一つ重要な手法である。このような手法は一定の数のモーメント条件しかない場合、基本的に容易に利用できるが、モーメントの数が数多く存在する場合、すべてのモーメント条件を利用することは殆どの場合推定量の性質を低下させるため、モーメント条件の選び方を考える必要がある。この問題に関しては、Morimune [1983]は一つ重要な研究となっている。そして、森棟公夫 [1985], Bekker [1994], Donald and Newey [2001], Kuersteiner and Okui [2009]などはさらにその問題に関して研究を進めてきた。そして、Kuersteiner and Okui [2009]はモーメントが複数存在するという問題の特殊ケースである操作変数の選択の問題にモデル平均の手法を応用した。Kuersteiner and Okui [2009]は二段階最小二乗法の第一段階で Hansen [2007]が提案したモデル平均の手法を利用し操作変数の選択を行い、最終的には最小の損失関数の値をもたらす推定量を導出した。

モデル平均の手法をボラティリティモデルの推定やValue-at-Riskの分析に応用した実証研究が幾つか存在する。Brownlees and Gallo [2008]はSmoothed-AIC, BIC, FIC¹の手法を用いてボラティリティモデルの変数選択を行った。また、Pesaran et al. [2009]はモデル平均の手法を利用してValue-at-Riskの分析を行い、モデル平均の推定結果に利用できる診断テストを考案した。

4.2 モデル平均による予測

これから日本の株価指数のデータを利用して、モデル平均の一つの適用例を示す。日本の日経225のデータに関して、幾つかの候補となるARCH(Autoregressive Conditional Heteroscedasticity)型モデルを推定して、株価指数の実現ボラティリティの予測を行う。

ARCH型モデルは株価のボラティリティを分析するための最もポピュラーなモデルのクラスのひとつである。Engle [1982]によって時系列データの分散不均一性を捉えるためのモデルとしてARCHモデルが考案されてから、データの特性を様々な側面から捉えるため、数多くのバリエティが提案されてきた。ここでは、その中の幾つかの代表的なモデルをモデル平均の候補モデルとする。以下では、簡単にそれらの候補モデルを紹介する。

¹FICはFocused Information Criterionの略である。推定対象となる異なった母数の特質に配慮した情報量基準である。詳細に関してはHjort and Claeskens [2003]またはClaeskens and Hjort [2008]を参照されたい。

4.2.1 ARCH型モデル

Engle [1982]は時系列データの分散不均一性を捉えるためのモデルとして, ARCH(q)モデルを提唱した。ARCH(q)モデルにおいては, 条件付分散が不均一であることが仮定されている。 $t - 1$ 期までの情報集合 ψ_{t-1} に含まれる内生変数のラグと外生変数の線形結合 $\mathbf{x}_t \boldsymbol{\beta}$ によって確率変数 r_t の条件付平均があらわされる。また, その条件付分散を h_t とする。そのモデル式は

$$r_t = \mathbf{x}_t \boldsymbol{\beta} + \sqrt{h_t} z_t \quad \text{or} \quad r_t | \psi_{t-1} \sim N[\mathbf{x}_t \boldsymbol{\beta}, h_t] \quad (37)$$

$$h_t = \alpha_0 + \sum_{i=1}^q \alpha_i \varepsilon_{t-i}^2 \quad (38)$$

$$\varepsilon_t = \sqrt{h_t} z_t \quad z_t \sim N[0, 1] \text{ i.i.d..} \quad (39)$$

となる。モデル式の中の h_t , 場合によっては $\sqrt{h_t}$ は, ファイナンスモデルにおけるボラティリティに該当する。モデル式を見れば分かるように, ARCH(q)モデルはボラティリティが過去の予期せぬリターンのショックの二乗によって影響されることを表現している。 h_t は条件付分散なので, 非負でなければならない。従って, $\alpha_0 > 0, \alpha_i \geq 0$ という制約が必要となる。

実務上, ARCH(q)モデルを金融データに適用して, ボラティリティに対するショックの持続性を表現しようとすると, ARCH(q)モデルの次数が長くなるという問題が生じる。この問題を回避するために提案されたのはGARCH(p,q)モデル(Bollerslev [1986])である。GARCH(p,q)モデルはARCH(q)モデルにボラティリティのラグの項を追加したものである。GARCH(p,q)モデルのボラティリティに関する式は次のようになる。

$$h_t = \alpha_0 + \sum_{i=1}^q \alpha_i \varepsilon_{t-i}^2 + \sum_{j=1}^p \beta_j h_{t-j} \quad (40)$$

ARモデルにMA項を追加することで, ARモデルの次数を減らすことと同様に, GARCHモデルにおいてはボラティリティのラグ項を導入することで, 短い次数, 普通はGARCH(1,1)モデルで, データの特性を捉えることができるようになる。

GARCHモデルにおいても, ボラティリティの非負性を保つため, 各パラメーターに制約をおく必要がある。ARCHモデルにおける制約

$$\alpha_0 > 0, \alpha_i \geq 0 \quad i = 1, 2, \dots, q$$

に加えて, $\beta_j \geq 0, j = 1, 2, \dots, p$ の制約条件が必要となる。

実証研究を積み重ねるうちに、金融データのボラティリティに関する様々な特徴が明らかになってきた。その中の一つとして、ボラティリティ変動の非対称性があげられる。非対称性とは、前日の株価変動の方向の違いによって、ボラティリティの増減の大きさの具合が変わるという特性である。経験的には前日の株価が下落した場合には、前日の株価が上昇した場合よりも、今日のボラティリティは大きくなる傾向が強い。しかし、GARCHモデルでは、この非対称性を捉えることができない。そこで、GJR-GARCHモデル、EGARCHモデル、TGARCHモデルなどが提案された。

Nelson [1991]によって提案されたEGARCH(Exponential GARCH)モデルは非説明変数の h_t を $\log(h_t)$ に置き換えることによって、ボラティリティ変動の非対称性を捉えることを可能にし、またGARCHモデルにおけるパラメーターの非負制約を除去することにも成功した。EGARCH(p,q)モデルは以下の式で表される。

$$\log(h_t) = \alpha_0 + \sum_{i=1}^q \alpha_i [\theta z_{t-i} + \gamma(|z_{t-i}| - E|z_{t-i}|)] + \sum_{j=1}^p \beta_j \log(h_{t-j}) \quad (41)$$

$$\varepsilon_t = \sqrt{h_t} z_t \quad z_t \sim N[0,1] \text{ i.i.d.} \quad (42)$$

ここで、 $\theta < 0$ であれば、前述した経験的事実のボラティリティ変動の非対称性が表現されることになる。 $\ln(h_t)$ がマイナスの値をとることもできるので、パラメーターに関する非負制約が不要となる。ただし、EGARCHモデルには欠点もある。GARCHモデルと異なり、 ε_t の代わりに z_t でモデルを定式化したため、特定のショック ε_t のボラティリティの変動に与える影響がわかりにくくなるのである。

Liu and Morimune [2006]は株価の連続的な上昇または下落がボラティリティに与える影響を捉えるため、OGARCHモデルを提案した。OGARCH(1,1)モデルのボラティリティ式は

$$h_t = \alpha_0 + \alpha_1 \exp(\phi \gamma_{t-1}) \varepsilon_{t-1}^2 + \beta h_{t-1}, \quad (43)$$

となる。 $\alpha_0, \alpha_1, \beta$ と ϕ は未知パラメーターで、 $\alpha_0 \geq 0, \alpha_1 > 0, \beta > 0$ である。そして、 γ_{t-1} は第 t 期まで同じ符号のショックが連続的に現れた日数である。その定義は

$$\gamma_{t-1} \equiv i, \text{ if } \text{sign}(\varepsilon_{t-1}) = \cdots = \text{sign}(\varepsilon_{t-i}) = -\text{sign}(\varepsilon_{t-(i+1)}) \quad (44)$$

である。 γ_{t-1} は $1 \leq \gamma_{t-2} + 1$ という値をとる。モデル式の中の指數関数の項は非負で、係数 ϕ が正の値と期待される。すなわち、株価が同じ方向に連続的に動くと、ボラティリティが大きくなりがちであると期待する。

4.2.2 実現ボラティリティの予測

応用例では、GARCH(1,1), OGARCH(1,1)及びEGARCH(1,1)モデルを利用する。モデル平均はSmoothed BICを採用する。データは1984年1月9日から2007年12月25日までの日経225の日次データと週次データである。²日次と週次データのサンプルサイズはそれぞれ5902と1242である。日次データから週次の実現ボラティリティ (Realized Volatility, RV)を計算する。そして、週次データで上記ARCH型モデルを推定し、 RV を予測する。週次の株価を y_t , $t = 1, 2, \dots, n$ として、ログ・リターンは

$$r_t = 100 [\log(y_t) - \log(y_{t-1})] \quad (45)$$

とする。週次のRVを

$$RV_t = \sum_{i=1}^{n_t} r_{t_i}^2 \quad (46)$$

と定義する。ただし、 n_t は第 t 週間にあった営業日の数で、 r_{t_i} はその週の第 i 番目のログ・リターンとする。ローリング・ウインドの幅を442とする。442週間のデータを利用してモデルを推定し、一週間先の RV を予測する。ローリング・ウインドを移動しながら、残りの800週間の RV を予測する。予測の良さを計るため、平均予測誤差二乗和

$$MSPE = \frac{1}{n} \sum_{t=1}^n (\widehat{RV}_t - RV_t)^2 \quad (47)$$

を利用する。ただし、 \widehat{RV}_t は RV_t の推定値を表す。二通りのSmoothed BICによるモデル平均を行い、予測する。一つはGARCH(1,1)とOGARCH(1,1)との二つのモデルの平均でもう一つはGARCH(1,1), OGARCH(1,1)とEGARCH(1,1)との三つのモデルによる平均である。モデル平均による予測値は

$$\widetilde{RV}_t = \sum_{m=1}^M c_{BIC}(m) \widehat{RV}_{m,t} \quad (48)$$

と定義する。ただし、 M はモデル平均に使われるモデルの数で、 $c_{BIC}(m)$ はモデル m に対応するBICによって計算されるウェイトである。詳細は前述Smoothed BICに関する説明を参照のこと。結果はTable1にまとめた。

Table1. 予測の結果

	GARCH	OGARCH	EGARCH	Best-BIC-I	Best-BIC-II	2モデル	3モデル
MSPE	95.91	96.26	100.05	96.26	98.77	95.83	97.87

²データの出所：Morningstar, Inc..

Best-BIC-Iはローリング予測の中で、GARCH(1,1)とOGARCH(1,1)の中のBICが小さい方で予測した場合の結果を表す。Best-BIC-IIは三つのモデルの場合である。そして、2モデルと3モデルは二つのモデルと三つのモデルによるモデル平均をあらわす。単独のGARCHやOGRCHと比べて、単独のEGARCHのMSPEが一番大きい。そして、二つのモデルによるモデル平均の結果はBest-BIC-Iよりも良く一番いい結果を出している。三つのモデル平均の結果はEGARCHモデルのパフォーマンスに影響され悪い結果となっているが、それでもBest-BIC-IIよりよくなっている。

以上の結果は、モデル平均は単純なBICによるモデル選択よりもある程度優れていることを示しているが、実証研究を行う際、必ずしもいつも言えることではない。以上の例に関しても、ローリング・ウインドの幅を変えたり、予測期間を長くしたり、候補となるモデル群を変えたりすることで、異なる結果が出ることになる。モデル平均の長所は理論上のもので、それを実証研究に反映させるために、モデル平均以外の更なる工夫が必要となる。

5 結論

本稿ではモデル選択の拡張であるモデル平均に関して最近の研究成果を概観した。さらに、その未解決問題について論じて、応用例を挙げた。モデル平均はリスクの面で単純なモデル選択より優れている。モデル平均の漸近分布が導出され、モデル平均による推定量のリスクの評価の方法も確立されている。しかし、漸近分布が有限標本分布の良い近似とならず、統計的推論は困難である。この問題を解決すべく本稿は一つの方向を論文の中で述べた。この方向を含め、様々な試みが期待される。

参考文献

- Akaike, H. [1973] “Information theory and an extension of the maximum likelihood principle,” in N., P. B. and C. F. eds. *Proc. of the 2nd Int. Symp. on Information Theory*, pp. 267–281.
- Bekker, P. A. [1994] “Alternative Approximations to the Distributions of Instrumental Variable Estimators,” *Econometrica*, Vol. 62, No. 3, pp. 657–81, May.

- Beran, R. [2000] “REACT Scatterplot Smoothers: Superefficiency through Basis Economy,” *Journal of the American Statistical Association*, Vol. 95, No. 449, pp. 155-171.
- Beran, R. and L. Dumbgen [1998] “Modulation of estimators and confidence sets,” *Annals of Statistics*, Vol. 26, No. 5, pp. 1826-1856.
- Bollerslev, T. [1986] “Generalized autoregressive conditional heteroskedasticity,” *Journal of Econometrics*, Vol. 31, No. 3, pp. 307-327, April.
- Brownlees, C. T. and G. M. Gallo [2008] “On Variable Selection for Volatility Forecasting: The Role of Focused Selection Criteria,” *Journal of Financial Econometrics*, Vol. 6, No. 4, pp. 513-539.
- Buckland, S. T., C. Burnham, K. P. Burnham, and N. H. Augustin [1997] “Model selection: an integral part of inference,” *Biometrics*, Vol. 53, pp. 603–618.
- Claeskens, G. and N. L. Hjort [2008] *Model Selection and Model Averaging*: Cambridge University Press.
- Clyde, M. and E. I. George [2004] “Model Uncertainty,” *Statistical Science*, Vol. 19, No. 1, pp. 81-94.
- Donald, S. G. and W. K. Newey [2001] “Choosing the Number of Instruments,” *Econometrica*, Vol. 69, No. 5, pp. 1161-1191, September.
- Draper, D. [1995] “Assessment and Propagation of Model Uncertainty,” *Journal of the Royal Statistical Society. Series B (Methodological)*, Vol. 57, No. 1, pp. 45-97.
- Engle, R. F. [1982] “Autoregressive Conditional Heteroscedasticity with Estimates of the Variance of United Kingdom Inflation,” *Econometrica*, Vol. 50, No. 4, pp. 987-1007, July.
- Hansen, B. E. [2007] “Least Squares Model Averaging,” *Econometrica*, Vol. 75, No. 4, pp. 1175-1189, 07.
- Hjort, N. and G. Claeskens [2003] “Frequentist Model Average Estimators,” *Journal of the American Statistical Association*, Vol. 98, pp. 879-899, January.

- Hoeting, J. A., D. Madigan, A. E. Raftery, and C. T. Volinsky [1999] “Bayesian model averaging: a tutorial,” *Statistical Science*, Vol. 14, No. 4, pp. 382-417. with comments by M. Clyde, David Draper and E. I. George, and a rejoinder by the authors.
- Kuersteiner, G. and R. Okui [2009] “Instrument selection by first stage prediction averaging.” Unpublished manuscript.
- Kullback, S. and R. A. Leibler [1951] “On Information and Sufficiency,” *The Annals of Mathematical Statistics*, Vol. 22, No. 1, pp. 79–86.
- Leeb, H. and B. M. Potscher [2008] “Can One Estimate The Unconditional Distribution of Post-Model-Selection Estimators?” *Econometric Theory*, Vol. 24, No. 3, pp. 338-376.
- Liu, Q. and K. Morimune [2006] “A Modified GARCH Model with Spells of Shocks,” *Asia-Pacific Financial Markets*, Vol. 12, No. 1, pp. 29-44.
- Mallows, C. L. [1973] “Some comments on Cp,” *Technometrics*, Vol. 15, pp. 661–675.
- Morimune, K. [1983] “Approximate Distributions of k-Class Estimators When the Degree of Overidentifiability Is Large Compared with the Sample Size,” *Econometrica*, Vol. 51, No. 3, pp. 821-841, May.
- Nelson, D. B. [1991] “Conditional Heteroskedasticity in Asset Returns: A New Approach,” *Econometrica*, Vol. 59, No. 2, pp. 347-70, March.
- Pesaran, M. H., C. Schleicher, and P. Zaffaroni [2009] “Model averaging in risk management with an application to futures markets,” *Journal of Empirical Finance*, Vol. 16, No. 2, pp. 280–305.
- Potscher, B. M. [2006] “The Distribution of Model Averaging Estimators and an Impossibility Result Regarding Its Estimation,” MPRA Paper 73, University Library of Munich, Germany.
- Schwarz, G. [1978] “Estimating the Dimension of a Model,” *The Annals of Statistics*, Vol. 6, No. 2, pp. 461–464.

Stone, M. [1977] “An Asymptotic Equivalence of Choice of Model by Cross-Validation and Akaike’s Criterion,” *Journal of the Royal Statistical Society. Series B (Methodological)*, Vol. 39, No. 1, pp. 44–47.

小西貞則・北川源四郎 [2004] 『情報量基準』, 朝倉書店.

竹内啓 [1976] 「情報統計量の分布とモデルの適切さの基準」, 『数理科学』, 第153巻, 12-18頁.

森棟公夫 [1985] 『経済モデルの推定と検定』, 共立出版.