

LASSO における AIC の揺らぎとブートストラップ法による近似

統計数理研究所 モデリング研究系 坂田綾香

LASSO においては、AIC が予測誤差の不偏推定量となることが知られている。AIC はデータに応じて揺らぐため、モデル選択の信頼性を確かめるうえでも、揺らぎを定量的に把握することが重要である。このような動機下の研究は新しいものではないが、LASSO においては AIC の分布に関する性質が明らかにされていない。そこで① AIC の真の揺らぎの評価 ② ブートストラップ法による近似の精度評価について考える。線形モデルの LASSO は次のように定義される。

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

ここで $\mathbf{y} \in \mathbb{R}^M$ 、説明変数 $\mathbf{A} \in \mathbb{R}^{M \times N}$ 、スパース表現 $\mathbf{x} \in \mathbb{R}^N$ とし、 $M < N$ の場合を考える。

① AIC の真の揺らぎの評価

大偏差原理に従うと、AIC 値 $N\mathbf{a}$ の分布はレート関数 $R(\mathbf{a}) \leq 0$ により次のように与えられる。

$$P(\mathbf{a}) \propto \exp(NR(\mathbf{a}))$$

典型的な値は $R(\mathbf{a}) = 0$ に対応し、これは予測誤差に一致する。パラメータ n を導入して次のように母関数 $\phi(n)$ を定義すると、レート関数により次のように与えられる。

$$\phi(n) \equiv \frac{1}{N} \log \int d\mathbf{y} d\mathbf{A} \exp(-nNa(\mathbf{y}, \mathbf{A})) = \frac{1}{N} \log \int da \exp\{N(R(\mathbf{a}) - na)\} \rightarrow \max_{\mathbf{a}} R(\mathbf{a}) - na$$

ここで、積分を鞍点評価した。つまり、レート関数のルジャンドル変換が母関数となっているため、母関数に逆ルジャンドル変換を適用することで、レート関数を次のように導出できる。

$$R(\mathbf{a}) = \phi - n \frac{\partial \phi}{\partial n}$$

このように導出したレート関数を用いて、AIC のデータ揺らぎを明らかにする。

② ブートストラップ法による近似の精度評価

ここではノンパラメトリックなブートストラップ法を考える。サンプル $\{\mathbf{y}, \mathbf{A}\}$ の経験分布に従って生成されたブートストラップ標本を $\{\mathbf{y}^\dagger, \mathbf{A}^\dagger\}$ と表記する。ここで $\mathbf{y}^\dagger \in \mathbb{R}^{M_B}$ 、 $\mathbf{A}^\dagger \in \mathbb{R}^{M_B \times N}$ とする。全てのブートストラップ標本は、変数 $\mathbf{c} \in \{0, 1, \dots, M_B\}^M \mid \sum_{\mu} c_{\mu} = M_B$ を用いて表現できる。 c_{μ} は、データの μ 番目の成分がリサンプリングによって引かれる回数を意味する。AIC がある値 $N\mathbf{a}$ をとるブートストラップ標本の数のエントロピーを次のように定義する。

$$\omega(\mathbf{a}) \equiv \frac{1}{N} \log \sum_{\mathbf{c}} \delta(\mathbf{a}(\mathbf{c}) - \mathbf{a})$$

ブートストラップ標本に関する母関数を次のように定義すると、エントロピーはその逆ルジャンドル変換から導出できる。

$$\phi^\dagger(v) \equiv \frac{1}{N} \log \sum_{\mathbf{c}} \exp(-vNa) = \frac{1}{N} \log \int da \exp\{N(\omega(\mathbf{a}) - va)\} \rightarrow \max_{\mathbf{a}} \omega(\mathbf{a}) - va$$

これにより導出したエントロピーを用いて①のレート関数と比較し、ブートストラップ法による AIC の揺らぎの近似評価の精度を明らかにする。