

Model selection under SCAD and MCP based on approximate message passing

Institute of Statistical Mathematics Ayaka Sakata
Tokyo Institute of Technology Tomoyuki Obuchi

We study model selection for linear regression problems penalized by nonconvex penalties; SCAD and MCP. The problem is formulated as

$$\min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + J(\mathbf{x}; \lambda, a),$$

where $\mathbf{y} \in \mathbb{R}^M$ and $\mathbf{A} \in \mathbb{R}^{M \times N}$ denote response variables and predictor matrix, respectively. The regression coefficient to be estimated, $\mathbf{x} \in \mathbb{R}^N$, is penalized by $J(\mathbf{x}; \lambda, a)$, where λ and a are regularization parameters that control the form of SCAD and MCP. The density of non-zero components in the minimizer depends on the regularization parameters, and hence the determination of these parameters corresponds to the model selection. Hereafter, we denote the minimizer of the problem as $\hat{\mathbf{x}}(\mathbf{y}; \lambda, a)$.

The quantities to determine the regularization parameters considered here are in-sample error ϵ_{in} and extra-sample error ϵ_{ext} defined by

$$\epsilon_{in}(\mathbf{y}; \lambda, a) \equiv \frac{1}{M} E_{\mathbf{z}} \left[\sum_{\mu=1}^M (z_{\mu} - (\mathbf{A}\hat{\mathbf{x}}(\mathbf{y}; \lambda, a))_{\mu})^2 \right],$$

$$\epsilon_{ext}(\mathbf{y}; \lambda, a) \equiv E_{w, \tilde{\mathbf{A}}} \left[(w - \tilde{\mathbf{A}}\hat{\mathbf{x}}(\mathbf{y}; \lambda, a))^2 \right],$$

where the statistical property of \mathbf{z} in the in-sample error is equivalent to that of response variable \mathbf{y} . The generative process of variable w is also equivalent to that of the response variables, but the predictor vector $\tilde{\mathbf{A}}$, which is not contained in the training of $\hat{\mathbf{x}}(\mathbf{y}; \lambda, a)$, is assigned to w .

The regularization parameters that achieve the minimum of the in-sample or extra-sample error is appropriate for the description of the response variables. However, the exact computation of these quantities is impossible because we do not know the true generative process of the response variables. We propose numerical methods to construct the estimators of these quantities using approximate message passing.