

Learning with Noisy Labels: Variable Selection and Misclassification Probability Modeling

Hui GUO

Department of Computer Science
Western University, Canada

Joint work with Grace Y. Yi

Problem Formulation

Classification problem with label noise:

- Input vector: $X \in \mathbb{R}^p$
True label: $Y \in \{0, 1\}$
Noisy label: $Y^* \in \{0, 1\}$
- Main study (noisy) data: $\mathcal{D}_M \triangleq \{\{X_i, Y_i^*\} : i \in \mathcal{M}\}$ with size n
Validation data: $\mathcal{D}_V \triangleq \{\{X_i, Y_i, Y_i^*\} : i \in \mathcal{V}\}$ with size n_V
- *e.g., diagnoses using high-precision tools*
- **Goal:** Learning a valid classifier using \mathcal{D}_M and \mathcal{D}_V to accurately predict the true label for future inputs

Misclassification Probability Modeling

- **Parametric method:**

$$\text{logit } \gamma_{01}(x) = g_0(x; \nu); \text{ logit } \gamma_{10}(x) = g_1(x; \nu)$$

- $\theta = (\bar{\beta}^\top, \nu^\top)^\top$: all involved parameters, with true value $\theta_0 = (\bar{\beta}_0^\top, \nu_0^\top)^\top$
- Log-likelihood function: $\ell^* \Rightarrow \ell^*(\theta)$

- **Semiparametric method:**

$$\hat{\gamma}_{10}(x) = \frac{\sum_{i \in \mathcal{V}} y_i (1 - y_i^*) \tilde{K}_{i,x}}{\sum_{i \in \mathcal{V}} y_i \tilde{K}_{i,x}}; \quad \hat{\gamma}_{01}(x) = \frac{\sum_{i \in \mathcal{V}} (1 - y_i) y_i^* \tilde{K}_{i,x}}{\sum_{i \in \mathcal{V}} (1 - y_i) \tilde{K}_{i,x}}.$$

- $\tilde{K}_{i,x}$: kernel estimator

$$\tilde{K}_{i,x} = \underbrace{h^{-p_1} K\{(x_i^C - x^C)/h\}}_{\text{continuous components}} \cdot \underbrace{\omega^{\sum_{t=1}^{p_2} \mathbb{I}(x_{it}^D \neq x_t^D)}}_{\text{discrete components}}$$

- Log-likelihood function: $\ell^* \Rightarrow \hat{\ell}^*(\bar{\beta})$

Penalized Log-Likelihood Function Maximization

Penalized Log-Likelihood Function:

$$Q^* = \ell^* - n \sum_{j=1}^p \rho_{\lambda_n}(|\beta_j|)$$

$\ell^*(\theta)$ or $\hat{\ell}^*(\bar{\beta})$ ←

penalty function for variable selection

λ_n : nonnegative tuning parameter

e.g., SCAD, MCP

Theorem (Consistency).

Under regularity conditions, there exists a local maximizer of Q^* such that

$$\|\hat{\beta} - \bar{\beta}_0\|_2 = O_p(n^{-1/2} + a_n),$$

where $a_n = \max\{|\rho'_{\lambda_n}(|\beta_{j0}|)| : \beta_{j0} \neq 0, j = 1, \dots, p\}$.

Penalized Log-Likelihood Function Maximization

Theorem (Oracle Property).

Under regularity conditions, any \sqrt{n} -consistent estimator obtained from maximizing Q^* have the following properties: as $n \rightarrow \infty$,

- (i) (Sparsity) with probability tending to 1, $\widehat{\beta}_{\text{II}} = 0$,
- (ii) (Asymptotic normality)

- For the parametric method,

$$\sqrt{n}(l_{11}^* + \Sigma) \left\{ (\widehat{\beta}_I - \bar{\beta}_{I0}) + (l_{11}^* + \Sigma)^{-1} b \right\} \xrightarrow{d} \mathcal{N}(0, l_{11}^*).$$

- For the semiparametric method,

$$\sqrt{n} \{l_{11} + \Sigma\} \left[(\widehat{\beta}_I - \bar{\beta}_{I0}) + \{l_{11} + \Sigma\}^{-1} b \right] \xrightarrow{d} \mathcal{N}(0, l_{11} + \Sigma_{\text{sp}}).$$

Thank You!