



A Bachelor of Science thesis

Super Learners

and their oracle properties

Jinyang Liu

Supervised by Prof. Thomas Gerds
Co-supervised by Prof. Niels Richard Hansen
Department of Mathematical Sciences
University of Copenhagen, Denmark

Submitted: March 16, 2023

Contents

1	Introduction	3
2	The discrete super learner, dSL	6
2.1	Finite sample properties	6
3	The ensemble super learner, eSL	6
4	Simulation results	6
5	Discussion	6

1 Introduction

Our setup closely models what is described in [VDL06] and [LD03]. Let O_1, \dots, O_n be n -i.i.d. observations distributed according to $P \in \mathcal{P}$ on some measurable space $(\mathcal{O}, \mathcal{A})$ where $O_i \in \mathcal{O}$ for each i and \mathcal{P} is our statistical model. For a parameter set Θ we define the corresponding loss function $L : \mathcal{O} \times \Theta \rightarrow [0, \infty)$ as a measurable map such that our goal is to find an estimator $\hat{\theta}$ that minimizes the true risk function $R : \Theta \rightarrow \mathbb{R}$ given as

$$R(\theta) = \int L(x, \theta) dP(x) = EL(O_1)$$

The parameter set Θ can be Euclidean, but for the focus of this thesis we will consider it as a collection of functions of the form $\theta : \mathcal{O} \rightarrow \mathbb{R}$.

Example 1 (Regression functions Θ). Let $O_1 = (Y_1, X_1), \dots, O_n = (Y_n, X_n) \in \mathcal{O} = \mathbb{R} \times \mathcal{X}$ be i.i.d. observations such that satisfy the model

$$Y_1 = \theta_0(X_1) + \varepsilon,$$

for an unobservable stochastic error term ε . The goal is to estimate the **regression function** $\theta_0 \in \Theta$ where $\Theta = \{\theta : \mathcal{X} \rightarrow \mathbb{R}\}$, is the set of regression functions each having \mathcal{X} as their domain. [VDL06]

Example 2 (Parametric family). Consider the initial setup from example 1. If Y_i is $\mathcal{B}(\mathbb{R}) - \mathcal{B}(\mathbb{R})$ measurable and X_i is $\mathcal{F} - \mathcal{B}(\mathbb{R})$ measurable for some sigma-algebra \mathcal{F} on \mathcal{X} , then a **generalized regression model** could be considered as parametrized family of distributions given that Θ is finite-dimensional.

First parametrize the conditional distributions of Y_1 given $X_1 = x$ as $\mathcal{Q} = \{Q_{\theta(x), \eta} \mid \theta \in \Theta, \eta \in \mathbb{R}^m\}$ such that $Q_{\theta(x), \eta}$ is a valid probability distribution on $\mathcal{B}(\mathbb{R})$ for each $x \in \mathcal{X}$ and $(\theta, \eta) \in \Theta \times \mathbb{R}^m$. The parameter θ will be a parameter of interest, and η is some nuisance parameter. We have that

$$P(Y \in A \mid X = x) = Q_{\theta(x), \eta}(A) \quad \text{for all } A \in \mathcal{B}(\mathbb{R}).$$

Now assume that X_1 is distributed according to some H_0 on \mathcal{X} , then it is possible to identify a joint distribution $P_{\theta, \eta}$ over our observations for each $(\theta, \eta) \in \Theta \times \mathbb{R}^m$ such that

$$P_{\theta, \eta}(X \in A, Y \in B) = \int_A Q_{\theta(x), \eta}(B) dH_0(x)$$

for every $A \in \mathcal{F}$ and $B \in \mathcal{B}(\mathbb{R})$. Our parametric model is therefore $\mathcal{P} = \{P_{\theta, \eta} \mid \theta \in \Theta, \eta \in \mathbb{R}^m\}$. [Bic+93]

Example 3 (Logistic regression model). Let $O_1 = (Y_1, X_1), \dots, O_n = (Y_n, X_n) \in \mathcal{O} = \{0, 1\} \times \mathcal{X}$ be i.i.d. observations, where Y_i is binary and $\mathcal{X} \subseteq \mathbb{R}^k$. We would like to estimate the parameter function $\theta \in \Theta = \{\theta \mid \theta : \mathcal{X} \rightarrow [0, 1]\}$

$$\theta(x) = E(Y_1 \mid X_1 = x) = P(Y_1 = 1 \mid X_1 = x),$$

In logistic regression we assume that $\theta(x) = P(Y_1 = 1 \mid X_1 = x) = \text{expit}(\beta x)$ for some $\beta \in \mathbb{R}^k$, and then the goal becomes to estimate the k -dimensional parameter β , in this case the \mathbb{R}^k parameter β completely determines θ , so Θ is also k -dimensional. The conditional distributions of Y_1 given $X_1 = x$ are Bernoulli distributions and can be parametrized as $\mathcal{Q} = \{\text{Ber}(\text{expit}(\beta x)) \mid \beta \in \mathbb{R}^k\}$. Now from example 2 we know that the statistical model of our observations be parametrized through β , in particular we have

$$\begin{aligned} P_\beta(Y_1 = 1, X_1 \in A) &= \int_A Q_{\theta(x)}(\{1\}) dH_0(x) \\ &= \int_A \text{expit}(\beta x) dH_0(x) \end{aligned}$$

If H_0 has density f w.r.t. Lebesgue measure, we can write

$$P_\beta(Y_1 = 1, X_1 \in A) = \int_A \text{expit}(\beta x) f(x) dm(x)$$

Definition 1 (Estimator of θ_0). An estimator for $\theta_0 \in \Theta$ is a measurable map $\hat{\theta} : \mathcal{X}^n \rightarrow \Theta$.

Definition 2 (Prediction model).

We would like to consider a set estimators $\{\hat{\theta}_q(P_n) \mid 1 \leq q \leq p\}$, where we find $\hat{\theta}_{\hat{q}}(P_n)$, which denotes the estimator that minimizes R and \hat{q} may depend on the observations. In order to find \hat{q} we have to proceed via cross validation. In cross validation, we randomly split our data into a training set and a test set. Let $S = (S_1, \dots, S_n) \in \{0, 1\}^n$ independent of X_1, \dots, X_n such that $S_i = 0$ indicates that X_i should be in the training set and $S_i = 1$ indicates that X_i belongs to the test set. We can define the empirical distributions over these two subsets, $P_{n,S}^0$ and $P_{n,S}^1$ as

$$\begin{aligned} P_{n,S}^0 &= \frac{1}{n_0} \sum_{i:S_i=0} \delta_{X_i} \\ P_{n,S}^1 &= \frac{1}{1 - n_0} \sum_{i:S_i=1} \delta_{X_i} \end{aligned}$$

Where n_0 would be the number of S_i 's that are marked 0.

Definition 3 (True risk of q 'th estimator averaged over splits). Given the data $X \in \mathcal{X}^n$ and a set of estimators $\{\hat{\theta}_q \mid 1 \leq q \leq p\}$, $p \in \mathbb{N}$. The risks of these estimator averaged over the splits specified by some S is given as a function of q

$$q \mapsto E_S \int L(x, \hat{\theta}_q(P_{n,S}^0)) dP(x) = E_S R(\hat{\theta}_q(P_{n,S}^0))$$

Where P is the true distribution for our data X .

Definition 4 (Oracle selector). The oracle selector is a function $\tilde{q} : \mathcal{X}^n \rightarrow \{1, \dots, p\}$ which finds the estimator that minimizes the true risk given our data $X \in \mathcal{X}^n$.

$$\tilde{q}(X) = \arg \min_{1 \leq q \leq p} E_S R(\hat{\theta}_q(P_{n,S}^0))$$

Where $P_{n,s}^0$ is the empirical distribution over the training set of X as specified by some split-variable S .

In light of the above definitions, we will define the cross-validation risk and the cross-validation selector for our estimators

Definition 5 (Cross-validation risk of i 'th estimator averaged over splits). Given the data $X \in \mathcal{X}^n$ and a set of estimators $\{\hat{\theta}_q \mid 1 \leq q \leq p\}, p \in \mathbb{N}$. The cross-validation risks of these estimator averaged over the splits specified by some S is given as a function of q

$$q \mapsto E_S \int L(x, \hat{\theta}_q(P_{n,S}^0)) dP_{n,s}^1(x) = E_S \hat{R}(\hat{\theta}_q(P_{n,S}^0))$$

Where $P_{n,S}^1$ is the empirical distribution over the validation set for our data X . We write \hat{R} for empirical risk over the validation set.

Definition 6 (Cross-validation selector). The cross-validation selector is a function $\hat{q} : \mathcal{X}^n \rightarrow \{1, \dots, p\}$ which finds the estimator that minimizes the cross-validation risk given our data $X \in \mathcal{X}^n$.

$$\hat{q}(X) = \arg \min_{1 \leq q \leq p} E_S \hat{R}(\hat{\theta}_q(P_{n,S}^0))$$

Where \hat{R} is the empirical risk over the validation set and $P_{n,s}^0$ is the empirical distribution over the training set of X as specified by some split-variable S .

We are interested in the risk difference between the cross-validation selector and the oracle selector, we remark that the optimal risk is attained at the true value θ_0

$$R(\theta_0) = \int L(x, \theta_0) dP(x),$$

and clearly it is the case that $R(\theta_0) \leq R(\hat{\theta})$ for any estimator $\hat{\theta}$ of θ_0 . Given a set of estimators we define the centered conditional risk as the difference

$$\begin{aligned} \Delta_S(\hat{\theta}_{\hat{q}}, \theta_0) &= R(\hat{\theta}_{\hat{q}}(P_{n,S}^0)) - R(\theta_0) \\ &= E_S \int L(x, \hat{\theta}_{\hat{q}}(P_{n,S}^0)) - L(x, \theta_0) dP(x) \end{aligned}$$

The following result is due to [LD03]:

Theorem 7 (Asymptotic equality). *The cross validation selector \hat{q} performs asymptotically as well as the oracle selector \tilde{q} in the sense that*

$$\frac{\Delta_S(\hat{\theta}_{\hat{q}}, \theta_0)}{\Delta_S(\hat{\theta}_{\tilde{q}}, \theta_0)} \rightarrow 1 \quad \text{in probability for } n \rightarrow \infty$$

2 The discrete super learner, dSL

2.1 Finite sample properties

3 The ensemble super learner, eSL

4 Simulation results

5 Discussion

Pellentesque tincidunt sodales risus, vulputate iaculis odio dictum vitae. Ut ligula tortor, porta a consequat ac, commodo non risus. Nullam sagittis luctus pretium. Integer vel nibh at justo convallis imperdiet sit amet ut lorem. Sed in gravida turpis. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Sed in massa vitae ligula pellentesque feugiat vitae in risus. Cras iaculis tempus mi, sit amet viverra nulla viverra pellentesque.