

# Project: Cross-validation for model selection (rough draft)

Diego Urdapilleta de la Parra

March 5, 2025

## 1 Setup and preliminary results

Let  $n, p_n$  be positive integers and  $\mathcal{D}_n := \{(y_i, \mathbf{x}_i) : i \in [n]\}$  be a set of independent data points drawn from a distribution  $\mathbb{P}_{y, \mathbf{x}}$  for  $(y, \mathbf{x}) \in \mathbb{R}^{1+p_n}$ . We treat the  $\mathbf{x}_i$  as predictors of the outcome  $y_i$ , and we assume a linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

where  $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_n]^\top \in \mathbb{R}^{n \times p_n}$  is the design matrix,  $\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]^\top$ , and  $\mathbf{e}$  is a mean-zero random vector with  $\text{Cov}(\mathbf{e}) = \sigma_2 \mathbf{I}_n$ .

### Assumptions

We assume that the following are generally satisfied:

$$\begin{aligned} \text{H1 :} \quad & \liminf_{n \rightarrow \infty} \frac{1}{n} \|\mathbf{M}_\alpha \mathbf{X} \boldsymbol{\beta}\|^2 > 0 \text{ for all } \alpha \in \mathcal{A} \\ \text{H2 :} \quad & \mathbf{X}^\top \mathbf{X} = O(n) \quad \text{and} \quad (\mathbf{X}^\top \mathbf{X})^{-1} = O(n^{-1}) \\ \text{H3 :} \quad & \lim_{n \rightarrow \infty} \max_{i \leq n} h_{ii, \alpha} = 0 \text{ for all } \alpha \in \mathcal{A} \end{aligned}$$

We consider the following setup for model selection. Let  $\mathcal{A} \subset 2^{[p_n]}$  be a family of index sets representing candidate models. For  $\alpha \in \mathcal{A}$ , we denote by  $p_n(\alpha)$  the cardinality of  $\alpha$  and consider the model

$$m_\alpha(\mathbf{X}) = \mathbf{X}_\alpha \boldsymbol{\beta}_\alpha,$$

where  $\mathbf{X}_\alpha$  is the sub-matrix of  $\mathbf{X}$  containing only the columns indexed by  $\alpha$ , and  $\boldsymbol{\beta}_\alpha$  is the coefficient vector containing only the entries indexed by  $\alpha$  in  $\boldsymbol{\beta}$ .

1. We say  $\alpha \in \mathcal{A}$  is *correct* if  $\mathbb{E}[\mathbf{y} \mid \mathbf{X}] \stackrel{\text{a.s.}}{=} m_\alpha(\mathbf{X})$ , and we denote by  $\mathcal{A}_c$  the set of correct models in  $\mathcal{A}$
2. We say  $\alpha \in \mathcal{A}$  is *wrong* if it is not correct, and we denote by  $\mathcal{A}_w$  the set of wrong models in  $\mathcal{A}$

3. We say  $\mathcal{A}$  is *embedded* if there exists an enumeration  $\alpha_1, \alpha_2, \dots, \alpha_k$  of all elements in  $\mathcal{A}$  such that

$$\alpha_1 \subset \alpha_2 \subset \dots \subset \alpha_k.$$

### Definition 1.1

For  $\alpha \in \mathcal{A}$ , let  $\hat{\beta}_\alpha$  be the OLS estimator of  $\beta_\alpha$  and  $\hat{m}_\alpha(\mathbf{X}) := \mathbf{X}_\alpha \hat{\beta}_\alpha$ . We denote the average squared error of  $\hat{m}_\alpha$  by

$$L_n(\alpha) := \frac{1}{n} \|\mathbb{E}[\mathbf{y} \mid \mathbf{X}] - \hat{m}_\alpha(\mathbf{X})\|^2.$$

Additionally, we write  $R_n(\alpha) := \mathbb{E}[L_n(\alpha) \mid \mathbf{X}]$ .

### Proposition 1.1

If we assume a linear model  $\mathbf{y} = \mathbf{X}\beta + \mathbf{e}$ , then

$$L_n(\alpha) = \frac{1}{n} \|H_\alpha \mathbf{e}\|^2 + \frac{1}{n} \|M_\alpha \mathbf{X}\beta\|^2 \quad \text{and} \quad R_n(\alpha) = \frac{1}{n} \sigma^2 p_n(\alpha) + \frac{1}{n} \|M_\alpha \mathbf{X}\beta\|^2,$$

where  $H_\alpha = \mathbf{X}_\alpha (\mathbf{X}_\alpha^\top \mathbf{X}_\alpha)^{-1} \mathbf{X}_\alpha^\top$  and  $M_\alpha = I_n - H_\alpha$ .

*Proof.* First, we have that

$$\begin{aligned} \|\mathbb{E}[\mathbf{y} \mid \mathbf{X}] - \hat{m}_\alpha(\mathbf{X})\|^2 &= \|\mathbf{X}\beta - \mathbf{X}_\alpha \hat{\beta}_\alpha\|^2 \\ &= \|\mathbf{X}\beta - H_\alpha(\mathbf{X}\beta + \mathbf{e})\|^2 \\ &= \|M_\alpha \mathbf{X}\beta - H_\alpha \mathbf{e}\|^2. \end{aligned}$$

Notice that  $M_\alpha \mathbf{X}\beta$  and  $H_\alpha \mathbf{e}$  are orthogonal:

$$\mathbf{e}^\top H_\alpha M_\alpha \mathbf{X}\beta = \mathbf{e}^\top H_\alpha (I_n - H_\alpha) \mathbf{X}\beta = \mathbf{e}^\top H_\alpha \mathbf{X}\beta - \mathbf{e}^\top H_\alpha \mathbf{X}\beta = 0.$$

Hence, the first part follows from the Pythagorean theorem.

For the second part, we note that  $\mathbb{E}[\|H_\alpha \mathbf{e}\|^2 \mid \mathbf{X}] = \sigma^2 p_n(\alpha)$  by the “trace trick”, where  $p_n(\alpha)$  denotes the size of model  $\alpha$ .  $\square$

### Proposition 1.2

Suppose that the set of correct candidate models  $\mathcal{A}_c \subset \mathcal{A}$  is non-empty, and let  $\alpha^*$  be the smallest correct model in  $\mathcal{A}_c$ . Then,  $\alpha^*$  minimizes  $R_n(\alpha)$  over  $\alpha \in \mathcal{A}$ .

*Proof.* Let  $\alpha \in \mathcal{A}$  be arbitrary and suppose that  $\alpha \in \mathcal{A}_c$ . Then,  $\mathbf{X}_\alpha \beta_\alpha = \mathbf{X}\beta$  and  $p_n(\alpha^*) \leq$

$p_n(\alpha)$ . Thus,

$$\begin{aligned} R_n(\alpha) &= \frac{1}{n} \sigma^2 p_n(\alpha) + \frac{1}{n} \|M_\alpha \mathbf{X} \boldsymbol{\beta}\|^2 \\ &= \frac{1}{n} \sigma^2 p_n(\alpha) + \frac{1}{n} \underbrace{\|M_\alpha \mathbf{X}_\alpha \boldsymbol{\beta}_\alpha\|^2}_0 \\ &= \frac{1}{n} \sigma^2 p_n(\alpha) \geq \frac{1}{n} \sigma^2 p_n(\alpha^*) = R_n(\alpha^*). \end{aligned}$$

Now suppose that  $\alpha \in \mathcal{A}_w$ . If  $p_n(\alpha) \geq p_n(\alpha^*)$ , the result follows by assumption H1. If  $p_n(\alpha) < p_n(\alpha^*)$ , then ... **MISSING**.  $\square$

## 2 Leave-One-Out CV [1]

In this section, we assume that  $p(\alpha) := p_n(\alpha)$  is constant for each  $\alpha \in \mathcal{A}$ .

### Definition 2.1

The LOOCV estimator of  $R_n(\alpha)$  is

$$\hat{R}_n^{(1)}(\alpha) := \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \mathbf{x}_{i\alpha}^\top \hat{\boldsymbol{\beta}}_\alpha}{1 - h_{ii,\alpha}} \right)^2$$

### Lemma 2.1 (Shao, 1993)

$$\hat{R}_n^{(1)}(\alpha) = \begin{cases} R_n(\alpha) + \sigma^2 + o_{\mathbb{P}}(1) & \text{if } \alpha \in \mathcal{A}_w \\ \frac{1}{n} \|M_\alpha \mathbf{e}\|^2 + \frac{2}{n} \sigma^2 p(\alpha) + o_{\mathbb{P}}(n^{-1}) & \text{if } \alpha \in \mathcal{A}_c \end{cases}$$

*Proof.* Using the Taylor expansion of  $1/(1-x)^2 = 1 + 2x + O(x^2)$ , we have

$$\frac{1}{(1 - h_{ii,\alpha})^2} = 1 + 2h_{ii,\alpha} + O_{\mathbb{P}}(h_{ii,\alpha}^2).$$

Thus,

$$\hat{R}_n^{(1)}(\alpha) = \underbrace{\frac{1}{n} \sum_{i=1}^n \left( y_i - \mathbf{x}_{i\alpha}^\top \hat{\boldsymbol{\beta}}_\alpha \right)^2}_{\xi_{\alpha,n}} + \underbrace{\frac{1}{n} \sum_{i=1}^n \left( 2h_{ii,\alpha} + O_{\mathbb{P}}(h_{ii,\alpha}^2) \right) \left( y_i - \mathbf{x}_{i\alpha}^\top \hat{\boldsymbol{\beta}}_\alpha \right)^2}_{\zeta_{\alpha,n}} \quad (1)$$

Let  $\xi_{\alpha,n}$  and  $\zeta_{\alpha,n}$  denote the first and second terms in (1), respectively. Note that

$$\begin{aligned}\xi_{\alpha,n} &= \frac{1}{n} \|M_\alpha \mathbf{X} \boldsymbol{\beta} + M_\alpha \mathbf{e}\|^2 \\ &= \frac{1}{n} (\|M_\alpha \mathbf{e}\|^2 + \|M_\alpha \mathbf{X} \boldsymbol{\beta}\|^2 + 2\mathbf{e}^\top M_\alpha \mathbf{X} \boldsymbol{\beta})\end{aligned}\tag{2}$$

$$\begin{aligned}&= \frac{1}{n} \|\mathbf{e}\|^2 + \|M_\alpha \mathbf{X} \boldsymbol{\beta}\|^2 + \frac{1}{n} \|H_\alpha \mathbf{e}\|^2 + \frac{2}{n} \mathbf{e}^\top M_\alpha \mathbf{X} \boldsymbol{\beta} \\ &= \frac{1}{n} \|\mathbf{e}\|^2 + \frac{1}{n} \|M_\alpha \mathbf{X} \boldsymbol{\beta}\|^2 + o_{\mathbb{P}}(1).\end{aligned}\tag{3}$$

The equality at (3) follows from the fact that  $\mathbb{E} [\|H_\alpha \mathbf{e}\|^2 \mid \mathbf{X}] = \sigma^2 p(\alpha)$  and

$$\mathbb{E} [\mathbf{e}^\top M_\alpha \mathbf{X} \boldsymbol{\beta} \mid \mathbf{X}]^2 = \sigma^2 \|M_\alpha \mathbf{X} \boldsymbol{\beta}\|^2 = O_{\mathbb{P}}(n), \quad (?)$$

so that  $1/n \|H_\alpha \mathbf{e}\|^2 \rightarrow_{\mathbb{P}} 0$  and  $2/n (\mathbf{e}^\top M_\alpha \mathbf{X} \boldsymbol{\beta}) = O_{\mathbb{P}}(1)$ . (?)

Since  $0 < h_{ii,\alpha} < 1$ ,  $2h_{ii,\alpha} + O_{\mathbb{P}}(h_{ii,\alpha}^2) \leq O_{\mathbb{P}}(\max_i h_{ii,\alpha})$ . Thus,

$$\zeta_{\alpha,n} \leq O_{\mathbb{P}}\left(\max_i h_{ii,\alpha}\right) \left(\frac{1}{n} \sum_{i=1}^n \left(y_i - \mathbf{x}_{i\alpha}^\top \hat{\boldsymbol{\beta}}_\alpha\right)\right)\tag{4}$$

(3) and (4) imply the first case in the Lemma.

If  $\alpha \in \mathcal{A}^c$ , it is easy to see from (2) that  $\xi_{\alpha,n} = 1/n \|M_\alpha \mathbf{e}\|^2$ , Furthermore,

$$\zeta_{\alpha,n} = \frac{2}{n} \sigma^2 p(\alpha) + o_{\mathbb{P}}(1), \quad (?)$$

proving the second case. □

### Proposition 2.2 (Shao, 1993)

Let  $\hat{\alpha}^{(1)}$  be the model minimizing  $\hat{R}_n^{(1)}(\alpha)$ .

1. Under H1, H2, and H3,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\hat{\alpha}^{(1)} \in \mathcal{A}_w) = 0.$$

2. For  $\alpha \in \mathcal{A}_c$  with  $\alpha \neq \alpha^*$ ,

$$\mathbb{P}\left(\hat{R}_n^{(1)}(\alpha) \leq \hat{R}_n^{(1)}(\alpha^*)\right) = \mathbb{P}\left(2(p(\alpha) - p(\alpha^*))\sigma^2 < \mathbf{e}^\top (H_\alpha - H_{\alpha^*})\mathbf{e}\right) + o_{\mathbb{P}}(1).$$

In particular, if  $\mathbf{e} \sim \mathcal{N}(0_n, \sigma^2 I_n)$ ,

$$\mathbb{P}\left(\hat{R}_n^{(1)}(\alpha) \leq \hat{R}_n^{(1)}(\alpha^*)\right) = \mathbb{P}(2k < \chi^2(k)) + o_{\mathbb{P}}(1)$$

for  $k = p(\alpha) - p(\alpha^*)$ .

3. If  $p(\alpha^*) \neq p$ ,

$$\lim_{n \rightarrow \infty} \mathbb{P}(\hat{\alpha}^{(1)} = \alpha^*) \neq 1.$$

*Proof.*

1. **MISSING**

2. The first part follows from Lemma 2.1 by algebraic manipulation. The second part follows by noting that, if  $\mathbf{e} \sim \mathcal{N}(0_n, \sigma^2 I_n)$ , then

$$\frac{\mathbf{e}^\top}{\sigma} (H_\alpha - H_{\alpha^*}) \frac{\mathbf{e}}{\sigma} \sim \chi^2(\text{tr}(H_\alpha - H_{\alpha^*})).$$

3. If  $p(\alpha^*) = p$ , then  $\mathcal{A}_c = \{\alpha^*\}$ . It follows from 1. that  $\lim_{n \rightarrow \infty} \mathbb{P}(\hat{\alpha}^{(1)} = \alpha^*) = 1$ . Conversely, if  $\lim_{n \rightarrow \infty} \mathbb{P}(\hat{\alpha}^{(1)} = \alpha^*) = 1$  **MISSING**  $\square$

### Corollary 2.3

LOOCV is not consistent. In particular, overfits with non-vanishing probability.

## 3 Shao, 1997 [2]

### Definition 3.1

Let  $\hat{\alpha}_n$  be the model selected by minimizing some criterion  $\hat{R}_n$  over  $\mathcal{A}$ , and let  $\alpha_n^*$  denote the model minimizing  $R_n$  over  $\mathcal{A}$ . We say  $\hat{R}_n$  is *consistent* if

$$\mathbb{P}\{\hat{\alpha} = \alpha^*\} \rightarrow 1$$

as  $n \rightarrow \infty$ . We say that  $\hat{R}_n$  is *asymptotically loss efficient* if

$$\frac{R_n(\hat{\alpha})}{R_n(\alpha_n^*)} \rightarrow 1 \quad \text{a.s.}$$

### Proposition 3.1 (Shao, 1997)

Suppose H1,  $p_n/n \rightarrow 0$ , and that  $\mathcal{A}_c$  is non-empty for all but finitely many  $n$ .

1. If  $|\mathcal{A}_c| = 1$  for all but finitely many  $n$ , then consistency is equivalent to efficiency in the sense of Definition 3.1
2. If  $p_n(\alpha_n^*) \not\rightarrow_{\mathbb{P}} \infty$ , then consistency is equivalent to efficiency in the sense of Definition 3.1

## References

- [1] J. Shao, “Linear Model Selection by Cross-validation,” *Journal of the American Statistical Association*, vol. 88, pp. 486–494, June 1993.
- [2] J. Shao, “An Asymptotic Theory for Linear Model Selection,” *Statistica Sinica*, vol. 7, pp. 221–264, Apr. 1997.