

# Leave one out cross validation

Zheng Yuan

June 2019

## 1 Prediction Error

Consider a linear model

$$y = x^T \beta + e, \quad (1)$$

where  $y$  is a response variable,  $x$  is a predictor vector of dimension  $p$ ,  $x^T$  denotes the transpose of  $x$ ,  $\beta$  is a vector of  $p$  unknown parameters and  $e$  is a random error with mean 0 and variance  $\sigma^2$ . Since sometimes some of the components in  $\beta$  may be 0, a more compact expression would be

$$y = x_\alpha^T \beta_\alpha + e, \quad (2)$$

where  $\alpha$  is a subset of  $d_\alpha$  distinct positive integers that are less or equal than  $p$  and  $\beta_\alpha$  is the subset of  $\beta$  whose components are indexed by the integers in  $\alpha$ . Each subset  $\alpha$  corresponds to a certain model, denoted by  $M_\alpha$ . After knowing whether each component of  $\beta$  is 0 or not, all the models  $M_\alpha$  can be classified into two categories:

- Category I: At least one nonzero component of  $\beta$  is not in  $\beta_\alpha$
- Category II:  $\beta_\alpha$  contains all nonzero components of  $\beta$

Suppose that we now split data set into two parts,  $\{(y_i, x_i), i \in s\}$  and  $\{(y_i, x_i), i \in s^c\}$ , where  $s$  is a subset of  $\{1, 2, \dots, n\}$  containing  $n_\nu$  integers and  $s^c$  is its complements containing  $n_s$  integers, so  $n_\nu + n_s = n$ . Then model  $M_\alpha$  is fitted using  $\{(y_i, x_i), i \in s^c\}$  and prediction error is evaluated using the validation data  $\{(y_i, x_i), i \in s\}$ . The average squared prediction error is

$$n_\nu^{-1} \|y_s - \hat{y}_{\alpha, s^c}\|^2 = n_\nu^{-1} \|(I_{n_\nu} - Q_{\alpha, s})^{-1} (y_s - X_{\alpha, s} \hat{\beta}_\alpha)\|^2, \quad (3)$$

where  $\|u\|^2 = u^T u$  for a vector  $u$ ,  $y_s$  is the  $n_\nu$  vector consists of the components of  $y$  indexed by  $i \in s$ ,  $\hat{y}_{\alpha, s^c}$  is the prediction of  $y_s$  using  $\{(y_i, x_i), i \in s^c\}$  by least squared method under model  $M_\alpha$ ,  $Q_{\alpha, s} = X_{\alpha, s}^T (X_\alpha^T X_\alpha)^{-1} X_{\alpha, s}^T$  is the least squares estimator of  $\beta_\alpha$  using whole data set.

## 2 Probability of selecting unnecessary large models

### 2.1 Estimate prediction error by LOOCV

First, let us start from the  $CV(1)$  estimate of  $\Gamma_{\alpha,n}$ , the true prediction error. In fact, we can evaluate this with the help of projection matrix by a trick in algebra.

The  $CV(1)$  estimate of  $\Gamma_{\alpha,n}$  can be written as

$$\begin{aligned}
\hat{\Gamma}_{\alpha,n}^{CV} &= \frac{1}{n} \sum_i [(1 - w_{i\alpha})^{-1} (y_i - x_{i\alpha}^T \hat{\beta}_\alpha)]^2 \\
&= \frac{1}{n} \sum_i [1 + 2w_{i\alpha} + O(w_{i\alpha}^2)] (y_i - x_{i\alpha}^T \hat{\beta}_\alpha)^2 \quad (\text{By } (1 - w_{i\alpha})^{-2} = 1 + 2w_{i\alpha} + O(w_{i\alpha}^2)) \\
&= \frac{1}{n} \sum_i (y_i - x_{i\alpha}^T \hat{\beta}_\alpha)^2 + \frac{1}{n} \sum_i [2w_{i\alpha} + O(w_{i\alpha}^2)] (y_i - x_{i\alpha}^T \hat{\beta}_\alpha)^2 \\
&= \frac{1}{n} \sum_i r_{i\alpha}^2 + \frac{1}{n} \sum_i [2w_{i\alpha} + O(w_{i\alpha}^2)] r_{i\alpha}^2 \quad (\text{where } r_{i\alpha} = y_i - x_{i\alpha}^T \hat{\beta}_\alpha) \\
&= \xi_{\alpha,n} + \phi_{\alpha,n}.
\end{aligned} \tag{4}$$

Here, we denote  $\xi_{\alpha,n} = \frac{1}{n} \sum_i r_{i\alpha}^2$ ,  $\phi_{\alpha,n} = \frac{1}{n} \sum_i [2w_{i\alpha} + O(w_{i\alpha}^2)] r_{i\alpha}^2$  and  $w_{i\alpha}$  is the  $i$ th diagonal element of the projection matrix  $P_\alpha$ .

### 2.2 Decomposition of $\xi_{\alpha,n}$

$$\begin{aligned}
\xi_{\alpha,n} &= \frac{1}{n} \sum_i r_{i\alpha}^2 \\
&= \frac{1}{n} (y - P_\alpha)^T (y - P_\alpha) \\
&= \frac{1}{n} y^T (I_n - P_\alpha)^T (I_n - P_\alpha) y \\
&= \frac{1}{n} (X\beta + e)^T (I_n - P_\alpha)^T (I_n - P_\alpha) (X\beta + e) \\
&= \frac{1}{n} (X\beta + e)^T (I_n - P_\alpha) (X\beta + e) \\
&= \frac{1}{n} e^T (I_n - P_\alpha) e + \frac{1}{n} \beta^T X^T (I_n - P_\alpha) X \beta + 2n^{-1} e^T (I_n - P_\alpha) X \beta \\
&= \frac{1}{n} e^T (I_n - P_\alpha) e + \Delta_{\alpha,n} + 2n^{-1} e^T (I_n - P_\alpha) X \beta.
\end{aligned} \tag{5}$$

Here, we denote  $\frac{1}{n} \beta^T X^T (I_n - P_\alpha) X \beta$  as  $\Delta_{\alpha,n}$  and  $P_\alpha$  is the projection matrix (hat matrix) for  $M_\alpha$ , i.e.,  $P_\alpha = X_\alpha (X_\alpha^T X_\alpha)^{-1} X_\alpha^T$ .

Thus, we obtain the decomposition of  $\xi_{\alpha,n}$  into three parts:  $\frac{1}{n} e^T (I_n - P_\alpha) e$ ,  $\Delta_{\alpha,n}$  and  $2n^{-1} e^T (I_n - P_\alpha) X \beta$ .

### 2.3 $\xi_{\alpha,n}$ for $M_\alpha$ in Category II

For  $M_\alpha$  in Category II, we have  $X\beta = X_\alpha\beta_\alpha$ , since  $\beta_\alpha$  contains all nonzero components of  $\beta$  in this case.

Using this condition, now we have

$$\begin{aligned}\Delta_{\alpha,n} &= \frac{1}{n}\beta_\alpha^T X_\alpha^T (I_n - P_\alpha) X_\alpha \beta_\alpha \\ &= \frac{1}{n}\beta_\alpha^T X_\alpha^T X_\alpha \beta_\alpha - \frac{1}{n}\beta_\alpha^T X_\alpha^T P_\alpha X_\alpha \beta_\alpha,\end{aligned}\tag{6}$$

where

$$\begin{aligned}\frac{1}{n}\beta_\alpha^T X_\alpha^T P_\alpha X_\alpha \beta_\alpha &= \frac{1}{n}\beta_\alpha^T X_\alpha^T X_\alpha (X_\alpha^T X_\alpha)^{-1} X_\alpha^T X_\alpha \beta_\alpha \\ &= \frac{1}{n}\beta_\alpha^T X_\alpha^T X_\alpha \beta_\alpha\end{aligned}\tag{7}$$

Therefore,

$$\Delta_{\alpha,n} = 0.\tag{8}$$

Also,

$$\begin{aligned}\frac{1}{n}e^T (I_n - P_\alpha) X\beta &= \frac{1}{n}e^T (I_n - P_\alpha) X_\alpha \beta_\alpha \\ &= \frac{1}{n}e^T X_\alpha \beta_\alpha - \frac{1}{n}e^T P_\alpha X_\alpha \beta_\alpha \\ &= \frac{1}{n}e^T X_\alpha \beta_\alpha - \frac{1}{n}e^T X_\alpha (X_\alpha^T X_\alpha)^{-1} X_\alpha^T X_\alpha \beta_\alpha \\ &= \frac{1}{n}e^T X_\alpha \beta_\alpha - \frac{1}{n}e^T X_\alpha \beta_\alpha \\ &= 0\end{aligned}\tag{9}$$

Therefore we have

$$\xi_{\alpha,n} = \frac{1}{n}e^T (I_n - P_\alpha)e$$

for all the models  $M_\alpha$  in Category II.

### 2.4 Approximation for $\phi_{\alpha,n}$

From the notation above,

$$\begin{aligned}\phi_{\alpha,n} &= \frac{1}{n} \sum_i [2w_{i\alpha} + O(w_{i\alpha}^2)] r_{i\alpha}^2 \\ &= \frac{2}{n} \sum_i w_{i\alpha} r_{i\alpha}^2 + \frac{1}{n} \sum_i O(w_{i\alpha}^2) r_{i\alpha}^2\end{aligned}\tag{10}$$

Under the condition

$$\lim_{n \rightarrow \infty} \max_{i \leq n} w_{i\alpha} = 0,$$

we have

$$\frac{1}{n} \sum_i O(w_{i\alpha}^2) r_{i\alpha}^2 = o(n^{-1})$$

On the other hand,

$$\begin{aligned} \frac{2}{n} \sum_i w_{i\alpha} r_{i\alpha}^2 &= \frac{2}{n} \sum_i w_{i\alpha} (\sigma^2 + o(1)) \\ &= \frac{2\sigma^2}{n} \sum_i w_{i\alpha} + o(n^{-1}) \\ &= \frac{2\sigma^2}{n} \text{tr}(P_\alpha) + o(n^{-1}) \\ &= \frac{2\sigma^2}{n} \text{tr}(X_\alpha (X_\alpha^T X_\alpha)^{-1} X_\alpha^T) + o(n^{-1}) \\ &= \frac{2\sigma^2}{n} \text{tr}(X_\alpha^T X_\alpha (X_\alpha^T X_\alpha)^{-1}) + o(n^{-1}) \\ &= \frac{2\sigma^2}{n} \text{tr}(I_{d_\alpha}) + o(n^{-1}) \\ &= \frac{2}{n} d_\alpha \sigma^2 + o(n^{-1}). \end{aligned} \tag{11}$$

Therefore,

$$\phi_{\alpha,n} = \frac{2}{n} d_\alpha \sigma^2 + o(n^{-1}), \tag{12}$$

and we obtain an approximation for  $\phi_{\alpha,n}$

$$\phi_{\alpha,n} \approx \frac{2}{n} d_\alpha \sigma^2. \tag{13}$$

## 2.5 Two lemmas on projection matrix $P_\alpha$

**Lemma 1.**  $\alpha \in \mathcal{A}$  is a subset of  $d_\alpha$  distinct positive integers that are less or equal than  $p$ ,  $\beta_\alpha$  is the subset of  $\beta$  whose components are indexed by the integers in  $\alpha$  and similar as  $X_\alpha$ .  $M_\alpha$  is the model corresponds to  $\alpha$ . Let  $P_\alpha$  be the projection matrix for  $M_\alpha$ , i.e.,  $P_\alpha = X_\alpha (X_\alpha^T X_\alpha)^{-1} X_\alpha^T$ . Then the eigenvalues of  $P_\alpha$  can only be 0 or 1.

*Proof.* Let  $\lambda$  be the eigenvalue of  $P_\alpha$ . Then there exists  $v \in R^n$  such that

$$\lambda v = P_\alpha v$$

which implies

$$P_\alpha(\lambda v) = P_\alpha(P_\alpha v)$$

Note that

$$P_\alpha(P_\alpha v) = P_\alpha v = \lambda v$$

and

$$P_\alpha(\lambda v) = \lambda(P_\alpha v) = \lambda^2 v$$

which follows

$$\lambda v = \lambda^2 v.$$

By solving this,  $\lambda$  can only be 0 or 1.  $\square$

**Lemma 2.** *Under the notations in lemma 1, if  $e_0 \sim N(0, I_n)$  is a standard normal random vector, then  $e_0^T P_\alpha e_0 \sim \chi_{d_\alpha}^2$ , where  $d_\alpha$  is the number of elements in  $\alpha$ .*

*Proof.* First, let  $\lambda_1, \lambda_2, \dots, \lambda_n$  be the eigenvalues of  $P_\alpha$ . Using the proposition of eigenvalues, we have

$$\sum_i \lambda_i = \text{tr}(P_\alpha). \quad (14)$$

On the other hand, by the relation between matrix trace and rank, we have

$$\text{tr}(P_\alpha) = \text{rank}(P_\alpha) = d_\alpha. \quad (15)$$

Putting (14) and (15) together,

$$\sum_i \lambda_i = d_\alpha. \quad (16)$$

From lemma 1,  $\lambda_i$  are only be 0 or 1, so there are  $d_\alpha$  1's and  $n - d_\alpha$  0's in  $\lambda_i$ 's. Without loss of generality, we assume that the first  $d_\alpha$  eigenvalues equal to 1 and the others are 0. So we obtain the eigen decomposition for  $P_\alpha$ : there exists an  $R^{n \times n}$  matrix  $V$  satisfying  $V^T V = I_n$  and

$$P_\alpha = V^T D_\alpha V, \quad (17)$$

where  $D_\alpha$  is an  $R^{n \times n}$  diagonal matrix whose first  $d_\alpha$  elements on diagonal are 1 and the others are 0.

Now with this decomposition,

$$\begin{aligned} e_0^T P_\alpha e_0 &= e_0^T (V^T D_\alpha V) e_0 \\ &= (V e_0)^T D_\alpha (V e_0) \\ &= \sum_{i=1}^{d_\alpha} (V e_0)_{(i)}^2, \end{aligned} \quad (18)$$

where  $(V e_0)_{(i)}$  denotes the  $i$ th elements of vector  $V e_0$ .

Note that  $V e_0 \sim N(0, I_n)$ , so  $(V e_0)_{(i)}$ 's are independent standard normal random variables,  $i = 1, 2, \dots, d_\alpha$ . Therefore,

$$e_0^T P_\alpha e_0 = \sum_{i=1}^{d_\alpha} (V e_0)_{(i)}^2 \sim \chi_{d_\alpha}^2, \quad (19)$$

$\square$

## 2.6 Empirical probability of selecting an unnecessary complex model

After getting everything ready, now it's the time to put them together to prove the last theorem.

**Theorem 1.** *Suppose that*

$$\lim_{n \rightarrow \infty} \max_{i \leq n} w_{i\alpha} = 0$$

*holds, we have the following conclusions:*

(1) *If  $M_\alpha$  is in Category II, then*

$$\hat{\Gamma}_{\alpha,n} = \frac{1}{n} e^T e - \frac{1}{n} e^T P_\alpha e + \frac{2}{n} d_\alpha \sigma^2 + o(n^{-1}). \quad (20)$$

(2) *The empirical probability of selecting model  $M_\alpha$  is*

$$Pr(M_\alpha \text{ is preferable to } M_\star \text{ by the CV}(1)) = Pr(2k < \chi_k^2), \quad (21)$$

*where  $k$  is the difference between the dimension of  $M_\alpha$  and that of true model  $M_\star$ , namely,  $k = d_\alpha - d_\star$ ,  $\chi_k^2$  here denotes a chi-square random variable with  $k$  as its degree of freedom. Since we only consider the models in Category II,  $k \geq 0$ .*

*Proof.* After putting all the results in 2.1-2.4, the first conclusion can be derived immediately. Now let us focus on the second conclusion.

$$\begin{aligned} & Pr(M_\alpha \text{ is preferable to } M_\star \text{ by the CV}(1)) \\ &= Pr(\hat{\Gamma}_{\alpha,n}^{CV} < \hat{\Gamma}_{\star,n}^{CV}) \\ &= Pr(2n^{-1} d_\alpha \sigma^2 - n^{-1} e^T P_\alpha e < 2n^{-1} d_\star \sigma^2 - n^{-1} e^T P_\star e) \quad (\text{by conclusion (1)}) \\ &= Pr(2(d_\alpha - d_\star) \sigma^2 < e^T (P_\alpha - P_\star) e) \\ &= Pr(2(d_\alpha - d_\star) \sigma^2 < \sigma^2 e_0^T (P_\alpha - P_\star) e_0) \quad (\text{where } e = \sigma^2 e_0 \text{ and thus } e_0 \sim N(0, I_n)) \\ &= Pr(2(d_\alpha - d_\star) < e_0^T (P_\alpha - P_\star) e_0) \end{aligned} \quad (22)$$

Note that

$$e_0^T P_\alpha e_0 \sim \chi_{d_\alpha}^2$$

and

$$e_0^T P_\star e_0 \sim \chi_{d_\star}^2$$

which follows

$$e_0^T (P_\alpha - P_\star) e_0 \sim \chi_{d_\alpha - d_\star}^2 \quad (23)$$

since  $e_0^T (P_\alpha - P_\star) e_0$  is independent of  $e_0^T P_\star e_0$

Therefore,

$$\begin{aligned}
& Pr(M_\alpha \text{ is preferable to } M_\star \text{ by the CV}(1)) \\
&= Pr(2(d_\alpha - d_\star) < \chi_{d_\alpha - d_\star}^2) \\
&= Pr(2k < \chi_k^2).
\end{aligned} \tag{24}$$

□

## 2.7 A simple simulation study

One simple simulation verifying the theorem above is to just set the true model as the null model. In this case, we set response variable  $y \sim N(3, 1)$ , i.e.,

$$y = x_0 + e$$

, where  $x_0 = 3$  and  $e \sim N(0, \sigma^2)$ . Besides, we create a redundant variable  $x_1 \sim U(10, 20)$ .

In each time of simulation, with the distributions above, we can generate a data set of size 1000. Then we fit a null and a one-variable linear regression respectively using the data we generate. After this, we utilize leave one out cross validation to do the model selection.

It turns out that among 1000 times simulations, there are about 16.7% of them in which leave one out cross validation selects the unnecessary large model (the one-variable model) rather than the null model, which coincides with  $Pr(2 < \chi_2^2) = 0.157$ . (In this case,  $k=1$ ).