

# Linear Model Selection by Cross-Validation

JUN SHAO\*

We consider the problem of selecting a model having the best predictive ability among a class of linear models. The popular leave-one-out cross-validation method, which is asymptotically equivalent to many other model selection methods such as the Akaike information criterion (AIC), the  $C_p$ , and the bootstrap, is asymptotically inconsistent in the sense that the probability of selecting the model with the best predictive ability does not converge to 1 as the total number of observations  $n \rightarrow \infty$ . We show that the inconsistency of the leave-one-out cross-validation can be rectified by using a leave- $n_v$ -out cross-validation with  $n_v$ , the number of observations reserved for validation, satisfying  $n_v/n \rightarrow 1$  as  $n \rightarrow \infty$ . This is a somewhat shocking discovery, because  $n_v/n \rightarrow 1$  is totally opposite to the popular leave-one-out recipe in cross-validation. Motivations, justifications, and discussions of some practical aspects of the use of the leave- $n_v$ -out cross-validation method are provided, and results from a simulation study are presented.

KEY WORDS: Balanced incomplete; Consistency; Data splitting; Model assessment; Monte Carlo; Prediction.

## 1. INTRODUCTION

Cross-validation is a method for model selection according to the predictive ability of the models. Suppose that  $n$  data points are available for selecting a model from a class of models. The data set is split into two parts. The first part contains  $n_c$  data points used for fitting a model (model construction), whereas the second part contains  $n_v = n - n_c$  data points reserved for assessing the predictive ability of the model (model validation). Strictly speaking, model validation is carried out using not just  $n_v$ , but all the  $n = n_v + n_c$  data. There are  $\binom{n}{n_v}$  different ways to split the data set. Cross-validation, as its name indicates, selects the model with the best average predictive ability calculated based on all (or some) different ways of data splitting.

Clearly, the computational complexity of this method increases as  $n_v$  increases. That is why the simplest cross-validation with  $n_v = 1$  has been the main focus of researchers' attention over the past 30 years. Discussions and theoretical studies about the cross-validation method with  $n_v = 1$  under various situations can be found, for example, in Allen (1974), Stone (1974, 1977a,b), Geisser (1975), Wahba and Wold (1975), Efron (1983, 1986), Picard and Cook (1984), Herzberg and Tsukanov (1986), and Li (1987).

Throughout this article I assume that the number of predictors in each model under consideration does not change as  $n$  increases. In this case, it is known to many statisticians (although a rigorous statement has probably not been given in the literature) that the cross-validation with  $n_v = 1$  is asymptotically incorrect (inconsistent) and is too conservative in the sense that it tends to select an unnecessarily large model.

There are other methods for model selection, such as the Akaike information criterion (AIC) (Akaike 1974; Shibata 1981), the  $C_p$  (Mallows 1973), the jackknife, and the bootstrap (Efron 1983, 1986). All these methods are asymptotically equivalent to the cross-validation with  $n_v = 1$  (Stone 1977a; Efron 1983), however, and thus they share the same deficiency; that is, they are inconsistent.

In this article I show that in the problem of selecting linear models, this deficiency of the cross-validation with  $n_v = 1$  can be rectified by using a cross-validation with a large  $n_v$ .

(depending on  $n$ ). Our result is somewhat surprising; to have an asymptotically correct cross-validation procedure, we need to select  $n_v$  having the same rate of divergence to infinity as  $n$ ; that is,  $n_v/n \rightarrow 1$  as  $n \rightarrow \infty$ . The reason why such a large  $n_v$  is needed is explored, after taking a close look at the asymptotic behavior of the cross-validation procedures.

When  $n_v$  is large, the amount of computation required to use the cross-validation may be impractical. We consider a "balanced incomplete" cross-validation; that is, only a much smaller part of  $\binom{n}{n_v}$  splits are made according to a systematic manner. Two other alternatives—a Monte Carlo approximation and an analytic approximation to the leave- $n_v$ -out cross-validation—are also considered. Their performances are examined in a simulation study.

The issue of using more than one observation at a time in validation against leave-one-out was also raised by other researchers. Herzberg and Tsukanov (1986) did some simulation comparisons between the cross-validation procedures with  $n_v = 1$  and  $n_v = 2$ . They found that the leave-two-out cross-validation is sometimes better than the leave-one-out cross-validation, although the two procedures are asymptotically equivalent in theory. See also Geisser (1975), Burman (1989), and Zhang (1991). In the context of jackknife variance estimation for nonsmooth statistics (such as the sample quantiles), Shao and Wu (1989) showed that the inconsistency of the leave-one-out jackknife variance estimator can be rectified by using a leave- $n_v$ -out jackknife. The difference is that here we require that  $n_v/n \rightarrow 1$ , whereas in Shao and Wu (1989) the rate of  $n_v$  diverging to infinity was related to the smoothness of the given statistic.

It should be noted that the story is quite different in the cases where the number of predictors in one of the models under consideration increases as  $n$  increases. In such cases, Li (1987) showed that under some conditions, the leave-one-out cross-validation is consistent and is asymptotically optimal in some sense.

## 2. MODEL SELECTION AND PREDICTION ERROR

Consider a linear model

$$y = \mathbf{x}'\boldsymbol{\beta} + e, \quad (2.1)$$

\* Jun Shao is Associate Professor, Department of Mathematics, University of Ottawa, Ottawa K1N 6N5, Canada. The research was supported by Natural Science and Engineering Research Council of Canada. The author thanks the referees for thoughtful comments and helpful suggestions.















