Chapter 16 Selection of the Best Subset of Variables in Regression and Time Series Models

Nicholas A. Nechval University of Latvia, Latvia

Konstantin N. Nechval Transport and Telecommunication Institute, Latvia

> Maris Purgailis University of Latvia, Latvia

> Uldis Rozevskis University of Latvia, Latvia

ABSTRACT

The problem of variable selection is one of the most pervasive model selection problems in statistical applications. Often referred to as the problem of subset selection, it arises when one wants to model the relationship between a variable of interest and a subset of potential explanatory variables or predictors, but there is uncertainty about which subset to use. Several papers have dealt with various aspects of the problem but it appears that the typical regression user has not benefited appreciably. One reason for the lack of resolution of the problem, but rather several problems for which different answers might be appropriate. The intent of this chapter is not to give specific answers but merely to present a new simple multiplicative variable selection criterion based on the parametrically penalized residual sum of squares to address the subset of predictor variables without sacrificing any explanatory power. The variables, which optimize this criterion, are chosen to be the best variables. The authors find that the proposed criterion performs consistently well across a wide variety of variable selection problems. Practical utility of this criterion is demonstrated by numerical examples.

DOI: 10.4018/978-1-61520-668-1.ch016

INTRODUCTION

Variable selection refers to the problem of selecting input variables that are most predictive of a given outcome. Variable selection problems are found in all supervised or unsupervised machine learning tasks, classification, regression, time series prediction, pattern recognition.

In the recent years, variable selection has become the focus of considerable research in several areas of application for which datasets with tens or hundreds of thousands of variables are available. These areas include text processing, particularly in application to Internet documents, and genomics, particularly gene expression array data. The objective of variable selection is three-fold: to improve the prediction performance of the predictors, to provide faster and more cost-effective predictors, and to provide a better understanding of the underlying process that generated the data.

A number of studies in the statistical literature discuss the problem of selecting the best subset of predictor variables in regression. Such studies focus on subset selection methodologies, selection criteria, or a combination of both. The traditional selection methodologies can be enumerative (e.g. all subsets and best subsets procedures), sequential (e.g. forward selection, backward elimination, stepwise regression, and stagewise regression procedures), and screening-based (e.g. ridge regression and principal components analysis). Standard texts like Draper and Smith (1981) and Montgomery and Peck (1992) provide clear descriptions of these methodologies.

Some of the reasons for using only a subset of the available predictor variables (given by Miller, 2002) are:

- To estimate or predict at a lower cost by reducing the number of variables on which data are to be collected;
- To predict more accurately by eliminating uninformative variables;

- To describe multivariate data sets parsimoniously; and
- To estimate regression coefficients with smaller standard errors (particularly when some of the predictors are highly correlated).

These objectives are of course not completely compatible. Prediction is probably the most common objective, and here the range of values of the predictor variables for which predictions will be required is important. The subset of variables giving the best predictions in some sense, averaged over the region covered by the calibration data, may be very inferior to other subsets for extrapolation beyond this region. For prediction purposes, the regression coefficients are not the primary objective, and poorly estimated coefficients can sometimes yield acceptable predictions. On the other hand, if process control is the objective then it is of vital importance to know accurately how much change can be expected when one of the predictors changes or is changed.

Suppose that y, a variable of interest, and \mathbf{x}_1 , ..., \mathbf{x}_v , a set of potential explanatory variables or predictors, are vectors of *n* observations. The problem of variable selection, or subset selection as it is often called, arises when one wants to model the relationship between y and a subset of \mathbf{x}_1 , ..., \mathbf{x}_v , but there is uncertainty about which subset to use. Such a situation is particularly of interest when *v* is large and \mathbf{x}_1 , ..., \mathbf{x}_v is thought to contain many redundant or irrelevant variables.

The variable selection problem is most familiar in the linear regression context, where attention is restricted to normal linear models. Letting windex the subsets of $\mathbf{x}_1, ..., \mathbf{x}_v$ and letting p_w be the number of the parameters of the model based on the *w*th subset, the problem is to select and fit a model of the form

$$\mathbf{y} = \mathbf{X}_{w} \mathbf{a}_{w} + \boldsymbol{\varepsilon},\tag{1}$$

16 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/selection-best-subset-variables-

regression/39334

Related Content

The Dynamic Usage of Models (DYSAM) as a Theoretically-Based Phenomenological Tool for Managing Complexity and as a Research Framework

Gianfranco Minati (2010). Cybernetics and Systems Theory in Management: Tools, Views, and Advancements (pp. 176-190).

www.irma-international.org/chapter/dynamic-usage-models-dysam-theoretically/39328

Living with a Dam: A Case of Care Practices in Large Technical Systems

Tihomir Mitev (2015). International Journal of Actor-Network Theory and Technological Innovation (pp. 19-29).

www.irma-international.org/article/living-with-a-dam/128337

Have You Taken your Guys on the Journey?: An ANT Account of Information Systems Project Evaluation

Dubravka Cecez-Kecmanovicand Fouad Nagm (2009). International Journal of Actor-Network Theory and Technological Innovation (pp. 1-22).

www.irma-international.org/article/have-you-taken-your-guys/1375

Why Using Actor Network Theory (ANT) Can Help to Understand the Personally Controlled Electronic Health Record (PCEHR) in Australia

Imran Muhammad, Say Yen Teohand Nilmini Wickramasinghe (2012). *International Journal of Actor-Network Theory and Technological Innovation (pp. 44-60).* www.irma-international.org/article/using-actor-network-theory-ant/66877

Consortial Benchmarking: Applying an Innovative Industry-Academic Collaborative Case Study Approach in Systemic Management Research

Holger Schieleand Stefan Krummaker (2010). *Cybernetics and Systems Theory in Management: Tools, Views, and Advancements (pp. 93-107).*

www.irma-international.org/chapter/consortial-benchmarking-applying-innovative-industry/39324