

Penalized likelihood approach to variable selection for Cox's regression model under nested case-control sampling

Jie-Huei Wang^{1*} (王价辉), Chun-Hao Pan², I-Shou Chang³,
Chao Agnes Hsiung³ and Yi-Ching Wang⁴

¹ Institute of Statistics, National Tsing Hua University

² Department of Mathematics, National Central University

³ Population Health Sciences, National Health Research Institutes

⁴ Department of Pharmacology, National Cheng Kung University

Abstract

Assuming Cox's regression model, we consider penalized likelihood approaches to conduct variable selection under nested case-control sampling or case-cohort sampling, which are useful in the current study of genomic epidemiology. Penalized non-parametric maximum likelihood estimate (PNPMLE) are characterized by self-consistency equations derived from score functions, which form the basis of the algorithm to compute PNPML. A cross-validation method is used to choose the tuning parameter within a family of penalty function. Simulation studies indicate that the numerical performance of PNPML is satisfactory and that LASSO performs best when cohort size is small and SCAD performs best when cohort size is large and may eventually perform as well as the oracle estimator, resembling the findings when *i.i.d.* sampling is considered. This method is also illustrated to look for DNA methylation biomarkers. For SCAD, consistency, asymptotic normality and oracle properties of the PNPML, the sparsity property of the penalty, and a consistent estimate of the asymptotic variance, based on observed profile likelihood, are established.

Keywords: nested case-control sampling, oracle property, penalized maximum likelihood estimate, profile likelihood, SCAD, variable selection