

Abstract

We consider the problem of binary class probability estimation (CPE) when one class is rare compared to the other. It is well known that standard algorithms such as logistic regression do not perform well in this. Common fixes include under-sampling and weighting. Recently, Wang & Dey (2010) suggested the use of a parameterized family of asymmetric link functions based on the generalized extreme value (GEV) distribution. The approach showed promising initial results, but combined with the logarithmic CPE loss implicitly used in their work, it results in a non-convex loss that is difficult to optimize. In this paper, we use tools from the theory of proper composite losses (Buja et al., 2005; Reid & Williamson, 2010) to construct a canonical underlying CPE loss corresponding to the GEV link, which yields a convex proper composite loss that we call the GEV-canonical loss; this loss can be tailored for CPE when one class is rare, and can be easily minimized using an IRLS-type algorithm. Our experiments on both synthetic and real data demonstrate that the resulting algorithm outperforms common approaches such as under-sampling and weights-correction.



GEV-Canonical Regression for Accurate Binary Class Probability Estimation when One Class is Rare

Arpit Agarwal¹, Harikrishna Narasimhan¹, Shivaram Kalyanakrishnan² and Shivani Agarwal¹ ¹Indian Institute of Science, Bangalore ²Yahoo Labs, Bangalore

