

A Two-stage Subgroup Analysis

Chen Canyi

May 8, 2020

Table of contents

- 1 Introduction
- 2 Model
- 3 Algorithm
- 4 Emperical Example

Introduction

- Biological heterogeneity is common in many diseases.
- Thus there have been many attempts to identify subgroups from a heterogenous population.
- One of the most popular methods for finding subgroups from a heterogenous population is the mixture model.

The mixture model

- need specify the underlying distribution for data
- the number of subgroups
- not suitable for the high dimensional setting $p \gg n$, which is a common situation in precision medicine with high dimension genetic covariates p and relative small sample size n .

The model

- In this project, we consider the following model:

$$Y_i = \mu_i + \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, i = 1, \dots, n,$$

where Y_i is the response for the i -th subject, μ_i 's are unknown subject-specific intercepts.

- We assume the sparsity structure where only a subset \mathcal{S} of $\boldsymbol{\beta}$ is nonzero. We assume $|\mathcal{S}| = s$.

Historical literature

- In the situation $p < n$, Ma and Huang (2015) considered optimizing the following objective fucntion

$$Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}; \lambda) = 1/2 \sum_{i=1}^n (Y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \sum_{1 \leq i < j \leq n} p(|\mu_i - \mu_j|, \lambda).$$

- alternating direction method of multipliers (ADMM)
algorithm: a first order algorithm.

Historical literature (Cont.)

A nature extension to $p > n$

$$\tilde{Q}_n(\boldsymbol{\mu}, \boldsymbol{\beta}; \lambda) = 1/2 \sum_{i=1}^n (Y_i - \mu_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \alpha_1 \sum_{1 \leq i < j \leq n} p(|\mu_i - \mu_j|, \lambda) + \alpha_2 p(|\boldsymbol{\beta}|, \omega).$$

ADMM has to iterate on at least $\binom{n}{2} + p$ parameters sequentially.

The algorithm

- Instead we first use Lasso procedure to select related covariates from \mathbf{x}_i , then apply k-means algorithm to the adjusted \tilde{Y}_i . For the Lasso procedure, we minimize the following objective function

$$(\hat{\mu}(\lambda), \hat{\beta}(\lambda)) = \arg \min_{(\mu, \beta) \in \mathbb{R}^{p+1}} 1/2 \sum_{i=1}^n (Y_i - \mu - \mathbf{x}_i^T \beta)^2 + p(\beta, \lambda). \quad (1)$$

- $\tilde{Y}_i = Y_i - \mathbf{x}_i^T \hat{\beta}$

Choice of concave penalty function

- ℓ_1 -norm penalty (Tibshirani, 1996, LASSO), smoothly clipped absolute deviations penalty (Fan and Li, 2001, SCAD) or minimax concave penalty (Zhang, 2010, MCP).
- our choice: $p(\beta, \lambda) = \lambda|\beta|_1$.

Optimization methods

- The first order methods, such as Combettes and Pesquet (2011), Bach et al. (2012) and Tropp and Wright (2010), and Newton type algorithms, such as Fountoulakis et al. (2014) and Dassios et al. (2015), can be used to optimize (1).
- We suggest to use the primal dual active set method (Fan et al., 2014), which is in spirit a generalized version of Newton type method. It usually converges after one-step iteration if there is a very good initial value.

The remains: One-dimension cluster

- It remains to analysis the subgroups in $\tilde{\mathbf{Y}} = (\tilde{Y}_1, \dots, \tilde{Y}_n)^T$.
- e.g., k-means, k-median, k-modes, even the deep learning algorithm.
- It is NP-hard in a general Euclidean space, even when the number of clusters k is 2 (Aloise et al., 2009).
- optimal one-dimensional clustering: $O(qknp)$ v.s. $(kn \lg n)$

Selection of the tuning parameters k and λ

- We select the tuning parameter λ by validation.
- Bayesian Information Criterion

$$\text{BIC}(k) = \log \left[\sum_{i=1}^n \{Y_i - \hat{\mu}_i(k) - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}\}^2 / n \right] + C_n \log n / n(k+s),$$

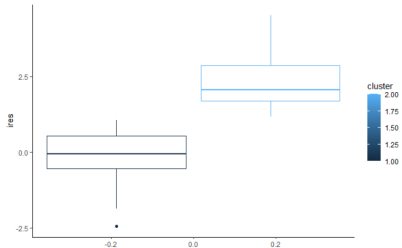
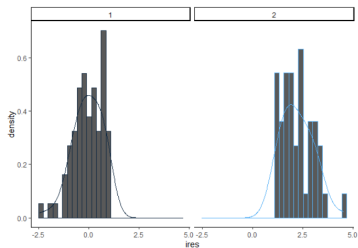
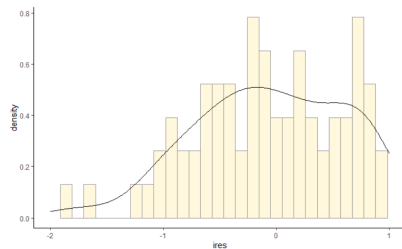
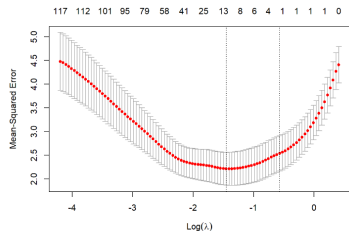
where $C_n = c \log\{\log(n+s)\}$ with some positive constant c .

We take $c = 10$ by defaults.

Chemores Data Example

- Lung cancer is one of the most prevalent and deadliest cancers, which can be classified into two major subtypes, small cell lung cancer (SCLC) and non-small cell lung cancer (NSCLC). NSCLC accouting 80lung cancers, is a known heterogeneous group and its prognosis is generally poor (Tsuboi et al., 2007).
- Publicly available lung cancer genomic data from the Chemores Cohort Study. This data is part of an integrated study of mRNA, miRNA and clinical variables to characterize the molecular distinctions between squamous cell carcinoma (SCC) and adenocarcinoma (AC) in Non Small Cell Lung Cancer (NSCLC) aside large cell lung carcinoma (LCC).

Results



Results

