

An algorithm for high-dimensional generalized linear mixed models using ℓ_1 -penalization

Jürg Schelldorfer

joint work with Peter Bühlmann

encouraged by Stephan Dlugosz

FOR916 Workshop on "Mathematical Statistics meets Econometrics"

June 9, 2011

Motivation

data set of Stephan Dlugosz:

administrative data set about employment:

Binary response variable $Y \in \{\text{employed, unemployed}\}$

covariates X : income, sex, age group, employment duration,....

quarterly results of (Y, X) of many workers over several years

General framework

- response variable from the exponential family
- continuous covariates
- grouped observations (think of longitudinal data, repeated measures data)

Goal:

Performing variable selection in the setup where AIC, BIC, cAIC, mAIC, ... are computationally infeasible (i.e. $n \approx p, n \ll p$)

Table of Contents

- 1 Introduction
- 2 High-dimensional Generalized Linear Mixed Models
- 3 Computational Algorithm
- 4 Outlook

Overview

	$n > p$	$n \ll p$
Generalized Linear Models (GLMs)	MLE [IRLS]	Lasso [R:glmnet]
Generalized Linear Mixed Models (GLMMs)	MLE [R:glmer]	?

n: number of observations

p: number of variables

Generalized Linear Model (GLM)

For n observations (y_i, x_i^T)

- (y_i, x_i^T) are independent for $i = 1, \dots, n$
- y_i has a density of the form

$$\exp \left\{ \phi^{-1} \left(y_i \xi_i - b(\xi_i) \right) + c(y_i, \phi) \right\} \text{ with } \mu_i = \mathbb{E}[y_i]$$

- $g(\mu) = \eta$ with $\eta = \mathbf{X}\beta$

Then estimate β by

$$\hat{\beta}_{MLE} = \operatorname{argmin}_{\beta} -\ell(\beta)$$

For $n \ll p$ we should not use the MLE. Use the Lasso (Tibshirani, 1996)

$$\hat{\beta}(\lambda) = \underset{\beta}{\operatorname{argmin}} -\ell(\beta) + \lambda \|\beta\|_1 \quad , \quad \lambda > 0$$

with the following properties:

- The Lasso does variable selection (i.e. some coefficients are set exactly to zero)
- Convex optimization problem, which can be solved efficiently

Generalized Linear Mixed Model (GLMM)

$g = 1, \dots, N$ independent groups/clusters/subjects

$j = 1, \dots, n_g$ observations for group/cluster/subject g

$n = \sum_{g=1}^N n_g$ total number of observations

\mathbf{y} : n -dim response variable

\mathbf{b} : q -dim (correlated) random effects

$\beta \in \mathbb{R}^p$ fixed-effects parameters

$\theta \in \mathbb{R}^L$ covariance parameters

ϕ dispersion parameter

\mathbf{X} : $n \times p$ model matrix for β

\mathbf{Z} : $n \times q$ model matrix for \mathbf{b}

Σ_θ : $q \times q$ covariance matrix, determined by θ

Generalized Linear Mixed Model (GLMM)

Model Assumptions:

- $y_i|\mathbf{b}$ are independent for $i = 1, \dots, n$
- $y_i|\mathbf{b}$ has a density of the form

$$\exp \left\{ \phi^{-1} \left(y_i \xi_i - b(\xi_i) \right) + c(y_i, \phi) \right\} \text{ with } \mu_i = \mathbb{E}[y_i|\mathbf{b}]$$

- $g(\mu) = \eta$ with $\eta = \mathbf{X}\beta + \mathbf{Z}\mathbf{b}$
- $\mathbf{b} \sim \mathcal{N}_q(\mathbf{0}, \Sigma_\theta)$ with $\Sigma_\theta \geq 0$ for $\theta \in \mathbb{R}^L$

$$(\hat{\beta}, \hat{\theta}, \hat{\phi})_{MLE} = \operatorname{argmin}_{\beta, \theta, \phi} -\log L(\beta, \theta, \phi)$$

Recap

	$n > p$	$n \ll p$
Generalized Linear Models (GLMs)	MLE [IRLS] ✓	Lasso [R:glmnet] ✓
Generalized Linear Mixed Models (GLMMs)	MLE [R:glmer] ✓	!

High-dimensional GLMM Set-up

Additionally to a GLMM, assume

- $n = \sum_{i=1}^N n_g \ll p$
- the true β_0 is sparse
- L small

Aim: Estimate β, θ, ϕ and predict \mathbf{b}

KeyIdea 1: Lasso-type penalty

Estimate the parameters (β, θ, ϕ) by minimizing

$$Q_\lambda(\beta, \theta, \phi) := -2 \log L(\beta, \theta, \phi) + \lambda \|\beta\|_1,$$

$$(\hat{\beta}, \hat{\theta}, \hat{\phi}) := \operatorname{argmin}_{\beta, \theta, \phi} Q_\lambda(\beta, \theta, \phi).$$

Remark: In general, $L(\beta, \theta, \phi)$ cannot be computed explicitly.

Laplace approximation

KeyIdea 2: Laplace approximation to approximate the integrand of $L(\beta, \theta, \phi)$ by a quadratic function.

$$I = \int_{\mathbb{R}^q} e^{-S(\mathbf{b})} d\mathbf{b} \approx (2\pi)^{q/2} |S''(\tilde{\mathbf{b}})|^{-1/2} e^{-S(\tilde{\mathbf{b}})}$$

where $\tilde{\mathbf{b}} = \operatorname{argmin}_{\mathbf{b}} S(\mathbf{b})$ is the mode of $-S(\mathbf{b})$.

Hence

$$Q_\lambda(\beta, \theta, \phi) \rightsquigarrow \tilde{Q}_\lambda^{LA}(\beta, \theta, \phi)$$

The GLMMLasso estimator

The GLMMLasso estimator is defined by

$$(\hat{\beta}, \hat{\theta}, \hat{\phi}) := \operatorname{argmin}_{\beta, \theta, \phi} \tilde{Q}_{\lambda}^{LA}(\beta, \theta, \phi)$$

Remark: It is a non-convex optimization problem!

How to calculate

$$(\hat{\beta}, \hat{\theta}, \hat{\phi}) := \operatorname{argmin}_{\beta, \theta, \phi} \tilde{Q}_{\lambda}^{LA}(\beta, \theta, \phi)?$$

KeyIdea 3: coordinate-wise optimization with inexact line search,

i.e. optimize \tilde{Q}_{λ}^{LA} w.r.t. one coordinate keeping all other coordinates fixed (Tseng and Yun, 2009):

- **Quadratic approximation** of the objective function
- calculate the **gradient**
- **Inexact line search** using the Armijo rule

The GLMMLasso algorithm II

GLMMLasso algorithm

(0) Choose a starting value $(\beta^{(0)}, \theta^{(0)}, \phi^{(0)})$.

Repeat for $s = 1, 2, \dots$

(1) (Fixed-effects parameter optimization) For $k = 1, \dots, p$

a) (Laplace approximation)

Calculate the Laplace approximation $\tilde{Q}_\lambda^{LA}(\cdot, \cdot, \cdot)$

b) (Quadratic approximation and inexact line search)

i) Approximate the second derivative by $h_k^{(s)} > 0$.

ii) Calculate the descent direction $d_k^{(s)} \in \mathbb{R}$

iii) Choose a step size $\alpha_k^{(s)} > 0$ such that there is a decrease in the objective function.

(2) (Covariance parameter optimization) For $l = 1, \dots, L$

$$\theta_l^{(s)} = \operatorname{argmin}_{\theta_l} \tilde{Q}_\lambda^{LA}(\cdot, \cdot, \cdot)$$

(3) (Dispersion parameter optimization)

$$\phi^{(s)} = \operatorname{argmin}_{\phi} \tilde{Q}_\lambda^{LA}(\cdot, \cdot, \cdot)$$

until convergence.

Tools to speed up

Two ingredients which speed up the algorithm remarkably:

- **KeyIdea 4a: regard $\tilde{\mathbf{b}}$ as fixed** for the quadratic approximation w.r.t. β_k
- **KeyIdea 4b: active-set algorithm** cycle through the non-zero coefficients β_k , and only through all p coefficients every D th iteration

This two ingredients make it feasible to calculate large data sets (i.e. $n = 400$ and $p = 4000$)!

The price to pay

Small additional bias in the parameter estimates, and similar variable selection properties.

This is ongoing work with Stephan Dlugosz on administrative data.

Take-home message

	$n > p$	$n \ll p$
Generalized Linear Models (GLMs)	MLE [IRLS]	Lasso [R:glmnet]
Generalized Linear Mixed Models (GLMMs)	MLE [R:glmer]	GLMMLasso KeyIdea 1-4

Thank you!

Questions?

References

- P. Tseng and S. Yun ; A Coordinate Gradient Descent Method for Nonsmooth Separable Minimization ; Mathematical Programming (2009)
- R. Tibshirani ; Regression Shrinkage and Selection via the Lasso ; J. R. Stat. Soc. (1996)
- J. Schelldorfer, P. Bühlmann and S. van de Geer ; Estimation for High-Dimensional Linear Mixed-Effects Models Using ℓ_1 -penalization ; The Scandinavian Journal of Statistics (2011)
- D. Bates ; lme4: Mixed-effects modeling with R (to appear)