

Stochastic Proximal Gradient Algorithm

Gersende FORT

LTCI, CNRS & Telecom ParisTech
Paris, France

Journée “Sélection dans les modèles mixtes”
Grenoble, Novembre 2014

Travail en collaboration avec Yves Atchadé (Univ. Michigan, USA) et Eric Moulines (Telecom ParisTech).

Plan

Régression Logistique à Effets Aléatoires

Le modèle

Vraisemblance pénalisée

Optimisation sous contrainte

Conclusion

Algorithme Gradient-Proximal

Algorithme Gradient-Proximal perturbé

Perturbation aléatoire type Monte Carlo

Régression Logistique à effets aléatoires (suite)

Conclusion

Régression logistique à effets aléatoires

- **Observations**: N observations binaires $\mathbf{Y} \in \{0,1\}^N$
- **Modèle à données cachées**: Conditionnellement à \mathbf{U} , pour tout $i = 1, \dots, N$,

$$Y_i \stackrel{\text{ind.}}{\sim} \mathcal{B} \left(\frac{\exp(\eta_i)}{1 + \exp(\eta_i)} \right)$$

où

$$\begin{bmatrix} \eta_1 \\ \dots \\ \eta_N \end{bmatrix} = \mathbf{X}\beta_{\text{true}} + \sigma_{\text{true}}\mathbf{Z}\mathbf{U}$$

- **Modèle paramétrique**:
 $\mathbf{X} \in \mathbb{R}^{N \times p}$, $\mathbf{Z} \in \mathbb{R}^{N \times q}$, connues déterministes.
 $\mathbf{U} \in \mathbb{R}^q$ vecteur aléatoire.
 $\beta_{\text{true}} \in \mathbb{R}^p, \sigma_{\text{true}} > 0$ paramètres inconnus.

Maximum de vraisemblance pénalisée

- **Log-vraisemblance** : Sous l'hypothèse $\mathbf{U} \sim \mathcal{N}_q(0, I)$, en posant

$$\theta = (\beta, \sigma) \quad F(\eta) = \frac{\exp(\eta)}{1 + \exp(\eta)}$$

la log-vraisemblance des observations \mathbf{Y} en θ est donnée par

$$\ell(\theta) = \log \int \prod_{i=1}^N \{F(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{ZU})_i)\}^{Y_i} \{1 - F(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{ZU})_i)\}^{1-Y_i} \phi(\mathbf{u}) d\mathbf{u}$$

- **Avec contraintes** : Dans un contexte de grande dimension $N \ll p$, introduction de pénalités

$$g_{\lambda, \theta}(\theta) = \lambda \left(\frac{1 - \alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$
$$\tilde{g}_{\mathcal{C}}(\theta) = \begin{cases} 0 & \text{si } \theta \in \mathcal{C} \\ +\infty & \text{sinon} \end{cases}$$

Minimisation sous contraintes (1/2)

$$\min_{\theta \in \Theta} (-\ell(\theta) + g(\theta))$$

Quelle régularité sur la log-vraisemblance $\theta \mapsto \ell(\theta)$?

- **Fonction :**

$$\ell(\theta) = \log \int \exp(\ell_c(\theta|\mathbf{u})) \phi(\mathbf{u}) d\mathbf{u}$$

avec

$$\ell_c(\theta|\mathbf{u}) = \sum_{i=1}^N \{Y_i (\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{u})_i) - \ln(1 + \exp(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{u})_i))\}$$

- **Gradient :**

$$\nabla \ell(\theta) = \int \left\{ \sum_{i=1}^N \{Y_i - F(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{u})_i)\} \begin{bmatrix} \mathbf{X}_i \cdot \\ (\mathbf{Z}\mathbf{u})_i \end{bmatrix} \right\} \pi_\theta(\mathbf{u}) d\mathbf{u}$$

avec

$$\pi_\theta(\mathbf{u}) = \exp(\ell_c(\theta|\mathbf{u}) - \ell(\theta)) \phi(\mathbf{u})$$

Minimisation sous contraintes (1/2)

$$\min_{\theta \in \Theta} (-\ell(\theta) + g(\theta))$$

Quelle régularité sur la log-vraisemblance $\theta \mapsto \ell(\theta)$?

- **Fonction :**

$$\ell(\theta) = \log \int \exp(\ell_c(\theta|\mathbf{u})) \phi(\mathbf{u}) d\mathbf{u}$$

avec

$$\ell_c(\theta|\mathbf{u}) = \sum_{i=1}^N \{Y_i (\mathbf{X}_i \cdot \beta + \sigma(\mathbf{ZU})_i) - \ln(1 + \exp(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{ZU})_i))\}$$

- **Hessien :** pas de “signe” particulier mais

$$\exists L, \forall \theta, \theta' \quad \|\nabla \ell(\theta) - \nabla \ell(\theta')\| \leq L \|\theta - \theta'\|$$

Minimisation sous contraintes (1/2)

$$\min_{\theta \in \Theta} (-\ell(\theta) + g(\theta))$$

Quelle régularité sur la log-vraisemblance $\theta \mapsto \ell(\theta)$?

- Fonction :

$$\ell(\theta) = \log \int \exp(\ell_c(\theta|\mathbf{u})) \phi(\mathbf{u}) d\mathbf{u}$$

avec

$$\ell_c(\theta|\mathbf{u}) = \sum_{i=1}^N \{Y_i (\mathbf{X}_i \cdot \beta + \sigma(\mathbf{ZU})_i) - \ln(1 + \exp(\mathbf{X}_i \cdot \beta + \sigma(\mathbf{ZU})_i))\}$$

↔ Fonction non explicite, définie comme une espérance

Fonction pas nécessairement convexe

Gradient non explicite, défini comme une espérance

Gradient lipschitzien

Minimisation sous contraintes (2/2)

$$\min_{\theta \in \Theta} (-\ell(\theta) + g(\theta))$$

Quelle régularité sur la fonction $\theta \mapsto g(\theta)$?

$$g_{\lambda, \theta}(\theta) = \lambda \left(\frac{1 - \alpha}{2} \|\beta\|_2^2 + \alpha \|\beta\|_1 \right)$$
$$\tilde{g}_{\mathbb{C}}(\theta) = \begin{cases} 0 & \text{si } \theta \in \mathbb{C} \\ +\infty & \text{sinon} \end{cases} \quad \mathbb{C} \text{ convexe fermé}$$

↔ Fonction convexe,

Peut valoir $+\infty$ (en certains points)

Semi-continue inférieurement, pas nécessairement dérivable.

Conclusion

- Chercher le maximum de vraisemblance pénalisé dans ce modèle de régression logistique à effet mixte nécessite la résolution d'un problème d'optimisation sous contraintes.
- Plus précisément, résoudre

$$\operatorname{argmin}_{\theta \in \Theta} (-\ell(\theta) + g(\theta))$$

dans un contexte

- $-\ell$ pas nécessairement convexe, de gradient Lipschitz

$$\|\nabla \ell(\theta) - \nabla \ell(\theta')\| \leq L\|\theta - \theta'\|$$

- $-\ell, \nabla \ell$ ne sont pas explicites mais sont définis par des intégrales.
- g est convexe, pouvant valoir $+\infty$, semi-continue inférieurement

Plan

Régression Logistique à Effets Aléatoires

Algorithme Gradient-Proximal

Contexte

Algorithme Gradient Proximal

Résultat de Convergence

Conclusion

Algorithme Gradient-Proximal perturbé

Perturbation aléatoire type Monte Carlo

Régression Logistique à effets aléatoires (suite)

Conclusion

Problème considéré

$$\text{(P)} \quad \min_{\theta \in \Theta} F(\theta) \quad F(\theta) = -\ell(\theta) + g(\theta),$$

sous les hypothèses

(A0) Θ espace euclidien de dimension finie, avec norme $\|\cdot\|$ et produit scalaire $\langle \cdot, \cdot \rangle$.

(A1) $g : \Theta \rightarrow (-\infty, +\infty]$ convexe, non identiquement égale à $+\infty$, et semi-continue inférieurement.

La fonction $\ell : \Theta \rightarrow \mathbb{R}$ est C^1 et il existe L tq pour tout $\theta, \theta' \in \Theta$,

$$\|\nabla\ell(\theta) - \nabla\ell(\theta')\| \leq L\|\theta - \theta'\|,$$

où $\nabla\ell$ désigne le gradient de ℓ .

Algorithme Gradient Proximal (1/3)

- Algorithme de type gradient pour des cas où la fonction objectif n'est pas différentiable.
- Peut se comprendre comme une approche “majoration-minimisation”.

De la relation conséquence du gradient lipschitzien

$$-\ell(\vartheta) \leq -\ell(\theta) - \langle \nabla \ell(\theta), \vartheta - \theta \rangle + \frac{L}{2} \|\theta - \vartheta\|^2$$

on déduit pour tout $\gamma \in (0, 1/L]$

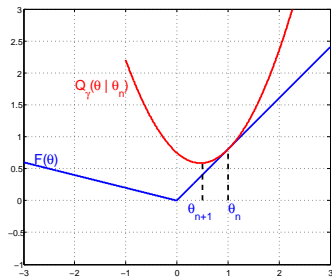
$$F(\vartheta) \leq -\ell(\theta) - \langle \nabla \ell(\theta), \vartheta - \theta \rangle + \frac{1}{2\gamma} \|\theta - \vartheta\|^2 + g(\vartheta)$$

Posons

$$\mathcal{Q}_\gamma(\vartheta|\theta) = -\ell(\theta) - \langle \nabla \ell(\theta), \vartheta - \theta \rangle + \frac{1}{2\gamma} \|\theta - \vartheta\|^2 + g(\vartheta)$$

Alors pour tout $\theta, \vartheta \in \Theta$ et $\gamma \in (0, 1/L]$

$$F(\vartheta) \leq \mathcal{Q}_\gamma(\vartheta|\theta) \qquad F(\theta) = \mathcal{Q}_\gamma(\theta|\theta)$$



$$F(\vartheta) \leq Q_\gamma(\vartheta|\theta_n)$$

$$F(\theta_n) = Q_\gamma(\theta_n|\theta_n)$$

Algorithme Gradient Proximal (2/3)

D'où l'algorithme

$$\begin{aligned}\theta_{n+1} &= \min_{\vartheta \in \Theta} Q_{\gamma}(\vartheta | \theta_n) \\ &= \min_{\vartheta \in \Theta} \left\{ g(\vartheta) + \frac{1}{2\gamma} \|\vartheta - (\theta_n + \gamma \nabla \ell(\theta_n))\|^2 \right\}\end{aligned}$$

où $\gamma \in (0, 1/L]$.

Algorithme Gradient Proximal:

$$\theta_{n+1} = \text{Prox}_{\gamma}(\theta_n + \gamma \nabla \ell(\theta_n))$$

avec

$$\text{Prox}_{\gamma}(\tau) = \min_{\theta \in \Theta} \left(g(\theta) + \frac{1}{2\gamma} \|\theta - \tau\|^2 \right)$$

Algorithme Gradient Proximal (2/3)

D'où l'algorithme

$$\begin{aligned}\theta_{n+1} &= \min_{\vartheta \in \Theta} Q_{\gamma}(\vartheta | \theta_n) \\ &= \min_{\vartheta \in \Theta} \left\{ g(\vartheta) + \frac{1}{2\gamma} \|\vartheta - (\theta_n + \gamma \nabla \ell(\theta_n))\|^2 \right\}\end{aligned}$$

où $\gamma \in (0, 1/L]$.

Algorithme Gradient Proximal: à pas décroissants

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1}}(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))$$

avec

$$\text{Prox}_{\gamma}(\tau) = \min_{\theta \in \Theta} \left(g(\theta) + \frac{1}{2\gamma} \|\theta - \tau\|^2 \right)$$

Algorithme Gradient Proximal (3/3)

Forme explicite du Prox dans certains cas. Par exemple,

- Lorsque $g(\theta) = 0$, \leftrightarrow gradient proximal = algo. de gradient

$$\theta_{n+1} = \theta_n + \gamma_{n+1} \nabla \ell(\theta_n)$$

Algorithme Gradient Proximal (3/3)

Forme explicite du Prox dans certains cas. Par exemple,

- Lorsque $g(\theta) = 0$, \leftrightarrow gradient proximal = algo. de gradient

$$\theta_{n+1} = \theta_n + \gamma_{n+1} \nabla \ell(\theta_n)$$

- Lorsque $g(\theta) = \begin{cases} 0 & \text{si } \theta \in C \\ +\infty & \text{sinon} \end{cases}$ pour un ensemble C fermé convexe,

$$\text{Prox}_\gamma(\tau) = \min_{\theta \in C} \|\tau - \theta\|^2 \quad \text{i.e. projection orthogonale sur } C$$

- \leftrightarrow gradient proximal = algo. de gradient projeté

$$\theta_{n+1} = \Pi_C(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))$$

Algorithme Gradient Proximal (3/3)

Forme explicite du Prox dans certains cas. Par exemple,

- Lorsque $g(\theta) = 0$, \hookrightarrow gradient proximal = algo. de gradient

$$\theta_{n+1} = \theta_n + \gamma_{n+1} \nabla \ell(\theta_n)$$

- Lorsque $g(\theta) = \begin{cases} 0 & \text{si } \theta \in C \\ +\infty & \text{sinon} \end{cases}$ pour un ensemble C fermé convexe,

$$\text{Prox}_\gamma(\tau) = \min_{\theta \in C} \|\tau - \theta\|^2 \quad \text{i.e. projection orthogonale sur } C$$

\hookrightarrow gradient proximal = algo. de gradient projeté

$$\theta_{n+1} = \Pi_C(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))$$

- Lorsque $g(\theta) = \lambda \left(\frac{1-\alpha}{2} \|\theta\|_2^2 + \alpha \|\theta\|_1 \right)$

$$(\text{Prox}_\gamma(\tau))_i = \frac{1}{1 + \gamma\lambda(1 - \alpha)} \begin{cases} \tau_i - \gamma\lambda\alpha & \text{si } \tau_i \geq \gamma\lambda\alpha \\ \tau_i + \gamma\lambda\alpha & \text{si } \tau_i \leq -\gamma\lambda\alpha \\ 0 & \text{sinon} \end{cases}$$

\hookrightarrow gradient proximal = algo. de gradient seuillé

$$\theta_{n+1} = \mathcal{S}_{\alpha, \lambda, \gamma_{n+1}}(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))$$

Convergence du Gradient-Proximal (1/4)

Quel ensemble limite \mathcal{L} ?

$$\theta_{n+1} = \text{Prox}_\gamma(\theta_n + \gamma \nabla \ell(\theta_n))$$

Lemma

Sous A0 et A1:

$$\mathcal{L} = \{\theta : \theta = \text{Prox}_\gamma(\theta + \gamma \nabla \ell(\theta))\} = \{\theta \in \text{Dom}(g) : 0 \in -\nabla \ell(\theta) + \partial g(\theta)\}.$$

Si de plus $-\ell$ est convexe

\mathcal{L} est l'ensemble des minimiseurs de F

Convergence du Gradient-Proximal (2/4)

Fonction de Lyapunov associée à un mapping continue.

- Par définition de la fonction $\mathcal{Q}_\gamma(\theta|\vartheta)$, on a :

$$F(\theta_{n+1}) \leq \mathcal{Q}_\gamma(\theta_{n+1}|\theta_n) \leq \mathcal{Q}_\gamma(\theta_n|\theta_n) = F(\theta_n)$$

- Régularité de $\theta \mapsto \text{Prox}_\gamma(\theta + \gamma \nabla \ell(\theta))$: il existe C tq

$$\|\text{Prox}_\gamma(\theta + \gamma \nabla \ell(\theta)) - \text{Prox}_\gamma(\theta' + \gamma \nabla \ell(\theta'))\| \leq C \|\theta - \theta'\|$$

Convergence du Gradient-Proximal (3/4)

Résultat général de convergence

Proposition (\dots ; AFM,14)

Sous A0-A1, pour $\gamma \in (0, 1/L]$. Si il existe \mathcal{K} compact tel que $\theta_n \in \mathcal{K}$ pour tout n , alors:

- (i) \mathcal{L} est non vide, et les valeurs d'adhérence de $\{\theta_n, n \geq 0\}$ sont dans $\mathcal{L} \cap \mathcal{K}$.
- (ii) Il existe $\theta_* \in \mathcal{L} \cap \mathcal{K}$ tel que $\lim_n F(\theta_n) = F(\theta_*)$.
- (iii) $\|\theta_{n+1} - \theta_n\| \rightarrow 0$.

- La suite peut rester dans un compact par construction (par ex. si g est la projection sur \mathcal{K})
- La suite reste dans un compact si $\{F \leq \theta_0\}$ est compact.
- Par (iii): soit $\{\theta_n, n \geq 0\}$ converge soit \mathcal{L} est un continuum.
- Par (ii,iii): $\{\theta_n\}$ converge dès que $\{\theta \in \mathcal{L} : F(\theta) = f_*\}$ est fini.

Convergence du Gradient-Proximal (4/4)

Résultat de convergence, dans le cas convexe

Proposition (\dots ; AFM,14)

Sous A0-A1 et $-\ell$ convexe; pour $\gamma \in (0, 1/L]$. Supposons une des conditions équivalentes

- (i) il existe un compact \mathcal{K} tel que $\theta_n \in \mathcal{K}$ pour tout n .*
- (ii) \mathcal{L} est non vide.*

Alors il existe $\theta_\star \in \mathcal{L}$ tel que $\lim_n \theta_n = \theta_\star$.

Conclusion

$$\text{(P)} \quad (\arg)\min_{\theta \in \Theta} \{-\ell(\theta) + g(\theta)\},$$

- L'algorithme Gradient-Proximal construit $\{\theta_n, n \geq 0\}$ tq sous certaines hypothèses
 - (cas convexe et non convexe) convergence de la fonction objectif $\{F(\theta_n), n \geq 0\}$
 - (cas convexe) converge vers θ_* , θ_* minimiseur de F .
 - (cas convexe) vitesses de convergence de $\{F(\theta_n), n \geq 0\}$ vers $F(\theta_*)$ en $1/n$.
- Sa mise en oeuvre nécessite de
 - bien choisir les pas γ_n
 - savoir calculer Prox_γ
 - calculer $\nabla \ell(\theta_n)$

↔ dans l'application statistique considérée au début, $\nabla \ell(\theta_n)$ n'est pas calculable. Quelle alternative?

Plan

Régression Logistique à Effets Aléatoires

Algorithme Gradient-Proximal

Algorithme Gradient-Proximal perturbé

Définition

Convergence de l'algorithme

Vitesses de convergence

Perturbation aléatoire type Monte Carlo

Régression Logistique à effets aléatoires (suite)

Conclusion

Algorithme Gradient-Proximal perturbé

- Algorithme exact :

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1}} (\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))$$

- Algorithme perturbé :

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1}} (\theta_n + \gamma_{n+1} H_{n+1})$$

où H_{n+1} est une approximation de $\nabla \ell(\theta_n)$.

Dans la suite,

- on étudie la suite $\{\theta_n\}$ produite par l'algorithme perturbé.
- On cherche des conditions suffisantes sur la perturbation $H_{n+1} - \nabla \ell(\theta_n)$ pour que l'algorithme perturbé converge.

Convergence de l'algorithme (1/2)

- La perturbation fait perdre la propriété de Lyapunov:

$$\begin{aligned} & F(\theta_{n+1}) - F(\theta_n) \\ &= F(\theta_{n+1}) - F(\text{Prox}_{\gamma_{n+1}}(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))) \\ &\quad + F(\text{Prox}_{\gamma_{n+1}}(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))) - F(\theta_n) \\ &= F(\text{Prox}_{\gamma_{n+1}}(\theta_n + \gamma_{n+1} H_{n+1})) - F(\text{Prox}_{\gamma_{n+1}}(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n))) + \underbrace{\dots}_{\text{quantité négative}} \end{aligned}$$

- Néanmoins, on sait montrer que sous A0-A1, pour tout compact \mathcal{K}

$$\lim_{n \rightarrow \infty} |F(\theta_{n+1}) - F(\text{Prox}_{\gamma_{n+1}}(\theta_n + \gamma_{n+1} \nabla \ell(\theta_n)))| \mathbb{1}_{\theta_n \in \mathcal{K}} = 0$$

dès que $\lim_n \{H_{n+1} - \nabla \ell(\theta_n)\} \mathbb{1}_{\theta_n \in \mathcal{K}} = 0$.

Convergence de l'algorithme (2/2)

Résultat de convergence, cas général

Théorème (AFM,14)

Sous A0-A1, pour $\gamma \in (0, 1/L]$.

Si \mathcal{L} est non vide, $\limsup_n \|\theta_n\|$ est finie et $\lim_n \{H_{n+1} - \nabla \ell(\theta_n)\} = 0$, alors la suite $\{F(\theta_n), n \geq 0\}$ converge vers une composante connexe de $F(\mathcal{L})$.

Si de plus $F(\mathcal{L})$ est d'intérieur vide, il existe $\theta_ \in \mathcal{L}$ tq*

(a) $\lim_n F(\theta_n) = F(\theta_*)$,

(b) la suite $\{\theta_n, n \geq 0\}$ converge vers $\mathcal{L} \cap \{\theta : F(\theta) = F(\theta_*)\}$.

- Quand $-\ell$ est convexe, $F(\mathcal{L})$ est d'intérieur vide. Si $-\ell$ fortement convexe, \mathcal{L} est un singleton.
- Dans le cas non convexe: résultat nouveau.
- Dans le cas convexe: conditions plus faibles que travaux antérieurs.

Vitesses de convergence (1/3)

Quelle stratégie? (a)

- Les méthodes Gdt-Prox proches de l'approximation stochastique

$$\begin{aligned}U_{n+1} &= U_n + \gamma_{n+1} H_{n+1} \\ &= U_n + \gamma_{n+1} \nabla h(U_n) + \gamma_{n+1} (H_{n+1} - \nabla h(U_n)).\end{aligned}$$

↔ perturbation non asympt. nulle et pas γ_n décroissant; techniques d'averaging améliorent la vitesse: $1/\sqrt{n}$ au lieu de $\sqrt{\gamma_n}$

Vitesses de convergence (1/3)

Quelle stratégie? (a)

- Les méthodes Gdt-Prox proches de l'approximation stochastique

$$\begin{aligned}U_{n+1} &= U_n + \gamma_{n+1} H_{n+1} \\ &= U_n + \gamma_{n+1} \nabla h(U_n) + \gamma_{n+1} (H_{n+1} - \nabla h(U_n)).\end{aligned}$$

↔ perturbation non asympt. nulle et pas γ_n décroissant; techniques d'averaging améliorent la vitesse: $1/\sqrt{n}$ au lieu de $\sqrt{\gamma_n}$

- Averaging pour Gdt-Prox : en parallèle, calculer récurivement

$$\bar{\theta}_n = \frac{\sum_{k=1}^n a_k \theta_k}{\sum_{k=1}^n a_k}$$

où $\{a_n, n \geq 0\}$ suite positive.

↔ Quelle suite $\{a_n, n \geq 0\}$ et a-t-on le même gain en vitesse que pour l'approx. sto.?

Vitesses de convergence (2/3)

Quelle stratégie? (b)

- Accélération “à la Nesterov”

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1}}(\vartheta_n + \gamma_{n+1} \nabla \ell(\vartheta_n))$$

où

$$\vartheta_n = \theta_n + \frac{(t_{n-1} - 1)^2}{t_n} (\theta_n - \theta_{n-1})$$

Vitesses de convergence (2/3)

Quelle stratégie? (b)

- Accélération “à la Nesterov”

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1}}(\vartheta_n + \gamma_{n+1} \nabla \ell(\vartheta_n))$$

où

$$\vartheta_n = \theta_n + \frac{(t_{n-1} - 1)^2}{t_n} (\theta_n - \theta_{n-1})$$

↪ pour l'algorithme exact, cas convexe : vitesse en $1/n^2$ au lieu de $1/n$

- Nesterov pour Gdt-Prox “perturbé”

$$\theta_{n+1} = \text{Prox}_{\gamma_{n+1}}(\vartheta_n + \gamma_{n+1} H_{n+1})$$

où H_{n+1} est une approximation de $\nabla \ell(\vartheta_n)$.

↪ A-t-on le même gain en vitesse que pour les algos Gdt-Prox non perturbés?

Vitesses de convergence (3/3)

Vitesse de convergence pour la suite averagée, cas convexe

Proposition (AFM,14)

Sous A0-A1 et $-\ell$ convexe, pour $\gamma_n \in (0, 1/L]$. Pour tout minimiseur θ_* de F ,

$$\begin{aligned} & \left(\sum_{j=1}^n a_j \right) (F(\bar{\theta}_n) - F(\theta_*)) \\ & \leq \frac{1}{2} \sum_{j=2}^n \left(\frac{a_j}{\gamma_j} - \frac{a_{j-1}}{\gamma_{j-1}} \right) \|\theta_{j-1} - \theta_*\|^2 + \frac{a_1}{2\gamma_1} \|\theta_1 - \theta_*\|^2 \\ & \quad + \sum_{j=1}^n a_j \langle \text{Prox}_{\gamma_j}(\theta_{j-1} + \gamma_j \nabla \ell(\theta_{j-1})) - \theta_*, \eta_j \rangle \\ & \quad + \sum_{j=1}^n a_j \gamma_j \|\eta_j\|^2. \end{aligned}$$

où $\eta_j = H_j - \nabla \ell(\theta_j)$.

Majoration analogue pour l'algo accéléré "à la Nesterov" voir AFM,14.

Plan

Régression Logistique à Effets Aléatoires

Algorithme Gradient-Proximal

Algorithme Gradient-Proximal perturbé

Perturbation aléatoire type Monte Carlo

Approximation Monte Carlo

Conditions Suffisantes pour la convergence

Vitesses de convergence

Régression Logistique à effets aléatoires (suite)

Conclusion

Approximation Monte Carlo (1/2)

Dans la suite, on considère le cas (plus général que l'apprentissage en ligne)

(A2)

$$\nabla \ell(\theta) = \int H_{\theta}(x) \pi_{\theta}(\mathbf{d}x)$$

Approximations possibles

- **Intégration numérique** - OK si espace de petite dimension
- **Echantillonnage d'importance**

$$\int H_{\theta}(x) \pi_{\theta}(x) \mathbf{d}x = \int H_{\theta}(x) \frac{\pi_{\theta}(x)}{\pi_{\star}(x)} \pi_{\star}(x) \mathbf{d}x \approx \frac{1}{m_{n+1}} \sum_{k=1}^{m_{n+1}} H_{\theta_n}(X_k) \frac{\pi_{\theta_n}(X_k)}{\pi_{\star}(X_k)}$$

- **MCMC**

$$\int H_{\theta}(x) \pi_{\theta}(x) \mathbf{d}x \approx \frac{1}{m_{n+1}} \sum_{k=1}^{m_{n+1}} H_{\theta_n}(X_{n,k})$$

où $\{X_{n,k}, k \geq 0\}$ est une chaîne de Markov de loi stationnaire $\pi_{\theta_n}(\mathbf{d}x)$.

Approximation Monte Carlo (2/2)

Dans cette partie, on considère l'approximation Monte Carlo

$$\nabla \ell(\theta) = \int H_{\theta}(x) \pi_{\theta}(x) dx \approx \frac{1}{m_{n+1}} \sum_{k=1}^{m_{n+1}} H_{\theta_n}(X_{n,k})$$

- Approximation biaisée

$$\mathbb{E}[H_{n+1} | \mathcal{F}_n] \neq \nabla \ell(\theta_n)$$

- mais pour certains échantillonneurs

$$|\mathbb{E}[H_{n+1} - \nabla \ell(\theta_n) | \mathcal{F}_n]| \leq \frac{C_{\theta_n}}{m_{n+1}} \quad \mathbb{E}[\|H_{n+1} - \nabla \ell(\theta_n)\|^2 | \mathcal{F}_n] \leq \frac{\tilde{C}_{\theta_n}}{m_{n+1}}$$

↔ Quels choix pour le pas γ_n , l'effort computationnel m_n : régime *approx. sto.* ou régime *mini-batch*?

Les versions stochastiques des algorithmes moyennés et à la Nesterov sont-elles aussi efficaces qu'en déterministe?

Convergence du Monte Carlo Gradient-Proximal

Les résultats

- de convergence
- de vitesse de convergence

énoncés précédemment restent valables : si les hypothèses sont vérifiées “presque-sûrement”, les conclusions sont vraies “presque-sûrement”.

En particulier, il existe des résultats dans la littérature sur la convergence des chaînes de Markov pour établir des conditions de la forme

$$\sum_n w_n (H_{n+1} - \nabla \ell(\theta_n)) \mathbb{1}_{\theta_n \in \mathcal{K}} < \infty \quad \text{p.s.}$$

pour tout compact \mathcal{K} et une suite déterministe $w_n \geq 0$.

Vitesse de convergence, cas convexe : algorithmes averagés (1/2)

Algorithme averagé

$$\bar{\theta}_n = \frac{\sum_{k=1}^n a_k \theta_k}{\sum_{k=1}^n a_k}$$

Proposition (AFM14)

Sous A0-A1-A2, $-\ell$ convexe et

$$\|\mathbb{E}[H_{n+1} - \nabla \ell(\theta_n) | \mathcal{F}_n]\| \leq \frac{C_1}{m_{n+1}} \quad \mathbb{E}[\|H_{n+1} - \nabla \ell(\theta_n)\|^2] \leq \frac{C_2}{m_{n+1}}$$

$$\sup_n \|\theta_n - \theta_\star\| \leq B \text{ p.s.}$$

Dans le cas $a_n = C_a n^a$ $a > -1$, $m_n = C_b n^b$ $b \geq 0$, $\gamma_n = C_c n^{-c}$ $c \geq 0$, pour tout $\theta_\star \in \mathcal{L}$,

$$\frac{\mathbb{E}[F(\bar{\theta}_n)] - F(\theta_\star)}{1 + a} \leq \frac{B^2}{2C_c} \frac{1}{n^{1-c}} + \frac{BC_1}{C_b} \frac{1}{n^{a+1}} \sum_{j=1}^n \frac{1}{j^{b-a}} + \frac{C_2 C_c}{C_b} \frac{1}{n^{a+1}} \sum_{j=1}^n \frac{1}{j^{b-a+c}}.$$

↔ quelle est la stratégie optimale pour (a,b,c) ?

Vitesse de convergence, cas convexe : algorithmes averagés (2/2)

	c	a	b	Rate $1/n^\bullet$	MC $1/\delta^\bullet$
Sans biais ($C_1 = 0$)	0	$(0, \infty)$	1	1	2
	$[0, 1/2]$	$(-c, +\infty)$	$1 - 2c$	$1 - c$	2
	$[0, 1)$	$-c$	$(1 - 2c, \infty) \cap [0, \infty)$	$1 - c$	$\frac{1+b}{1-c}$
Avec biais ($C_1 > 0$)	0	$(0, \infty)$	1	1	2
	$[0, 1)$	$(-c, \infty)$	$1 - c$	$1 - c$	$\frac{2-c}{1-c}$
	$[0, 1)$	$-c$	$(1 - c, \infty)$	$1 - c$	$\frac{1+b}{1-c}$
$C_1 = C_2 = 0$	$[0, 1)$	$(-1, \infty)$	-	$1 - c$	-

Valeurs de (a,b,c) pour atteindre la vitesse de convergence Rate.
 La colonne MC donne le nombre d'échantillons MC nécessaires pour avoir $\mathbb{E}[F(\bar{\theta}_n)] - F(\theta_*) = O(\delta)$. Comme référence, la dernière ligne donne la vitesse de l'algorithme exact (i.e. $H_{n+1} = \nabla \ell(\theta_n)$).

Vitesse de convergence, cas convexe : algorithmes averagés (2/2)

	c	a	b	Rate $1/n^\bullet$	MC $1/\delta^\bullet$
Sans biais ($C_1 = 0$)	0	$(0, \infty)$	1	1	2
	$[0, 1/2]$	$(-c, +\infty)$	$1 - 2c$	$1 - c$	2
	$[0, 1)$	$-c$	$(1 - 2c, \infty) \cap [0, \infty)$	$1 - c$	$\frac{1+b}{1-c}$
Avec biais ($C_1 > 0$)	0	$(0, \infty)$	1	1	2
	$[0, 1)$	$(-c, \infty)$	$1 - c$	$1 - c$	$\frac{2-c}{1-c}$
	$[0, 1)$	$-c$	$(1 - c, \infty)$	$1 - c$	$\frac{1+b}{1-c}$
$C_1 = C_2 = 0$	$[0, 1)$	$(-1, \infty)$	-	$1 - c$	-

Valeurs de (a,b,c) pour atteindre la vitesse de convergence Rate.
 La colonne MC donne le nombre d'échantillons MC nécessaires pour avoir $\mathbb{E}[F(\bar{\theta}_n)] - F(\theta_*) = O(\delta)$. Comme référence, la dernière ligne donne la vitesse de l'algorithme exact (i.e. $H_{n+1} = \nabla \ell(\theta_n)$).

Meilleure stratégie : en fonction du nombre d'itérations, en fonction du coût de calcul :

pas fixe $\gamma_n = \gamma$ croissance linéaire MC : $m_n \sim n$

Vitesse de convergence, cas convexe : accélération à la Nesterov (1/2)

Algorithme accéléré "à la Nesterov"

Proposition (AFM14)

Sous A0-A1-A2, $-\ell$ convexe et

$$|\mathbb{E}[H_{n+1} - \nabla\ell(\theta_n) | \mathcal{F}_n]| \leq \frac{C_1}{m_{n+1}} \quad \mathbb{E}[\|H_{n+1} - \nabla\ell(\theta_n)\|^2] \leq \frac{C_2}{m_{n+1}}$$

$$\sup_n \|\theta_n - \theta_\star\| \leq B \text{ p.s.}$$

Dans le cas $m_n = C_b n^b$ $b \geq 0$, $\gamma_n = C_c n^{-c}$ $c \geq 0$, pour tout $\theta_\star \in \mathcal{L}$,

$$\begin{aligned} \mathbb{E}[F(\theta_{n+1})] - F(\theta_\star) &\leq \frac{(n+1)^c}{t_n^2} \left(\mathbb{E}[F(\theta_1)] - F(\theta_\star) + \frac{B^2}{2C_c} \right) \\ &+ 4BC_1C_b \frac{(n+1)^c}{t_n^2} \sum_{k=1}^n \frac{t_k^2}{(k+1)^{c+b}} + \frac{C_2(n+1)^c}{t_n^2} \sum_{k=1}^n \frac{t_k^2}{(k+1)^{2c+b}}. \end{aligned}$$

↔ quelle est la stratégie optimale pour (b,c) ?

Vitesse de convergence, cas convexe : accélération à la Nesterov (2/2)

	c	b	Rate $1/n^\bullet$	MC $1/\delta^\bullet$
Sans biais ($C_1 = 0$)	0 [0,2)	(3,∞) (3 - 2c, +∞) ∩ [0,∞)	2 2 - c	(b + 1)/2 $\frac{b+1}{2-c}$
Avec biais ($C_1 > 0$)	0 [0,2)	(3,∞) (3 - c,∞)	2 2 - c	(b + 1)/2 $\frac{b+1}{2-c}$
$C_1 = C_2 = 0$	[0,2)	-	2 - c	-

Valeurs (b,c) pour atteindre la vitesse de convergence Rate quand $t_n = O(n^2)$. La colonne MC donne le nombre d'échantillons MC pour atteindre la précision $\mathbb{E}[F(\theta_n)] - F(\theta_*) = O(\delta)$. La dernière ligne donne, pour référence, la vitesse de l'algorithme exact (i.e. $H_{n+1} = \nabla \ell(\theta_n)$).

Vitesse de convergence, cas convexe : accélération à la Nesterov (2/2)

	c	b	Rate $1/n^\bullet$	MC $1/\delta^\bullet$
Sans biais ($C_1 = 0$)	0 [0,2)	(3,∞) $(3 - 2c, +\infty) \cap [0, \infty)$	2 $2 - c$	$(b + 1)/2$ $\frac{b+1}{2-c}$
Avec biais ($C_1 > 0$)	0 [0,2)	(3,∞) $(3 - c, \infty)$	2 $2 - c$	$(b + 1)/2$ $\frac{b+1}{2-c}$
$C_1 = C_2 = 0$	[0,2)	-	$2 - c$	-

Valeurs (b,c) pour atteindre la vitesse de convergence Rate quand $t_n = O(n^2)$. La colonne MC donne le nombre d'échantillons MC pour atteindre la précision $\mathbb{E}[F(\theta_n)] - F(\theta_*) = O(\delta)$. La dernière ligne donne, pour référence, la vitesse de l'algorithme exact (i.e. $H_{n+1} = \nabla \ell(\theta_n)$).

Meilleure stratégie : en fonction du nombre d'itérations, en fonction du coût de calcul :

pas fixe $\gamma_n = \gamma$ croissance MC : $m_n \sim n^{3+\iota}$

Vitesse de convergence, cas convexe : averaging ou Nesterov ?

- On peut choisir (a,b,c) pour que les versions stochastiques atteignent les mêmes vitesses que les algorithmes exacts.
- En particulier, vitesse maximale atteinte par
Averaging : $1/n$
Nesterov : $1/n^2$
que l'approximation soit avec ou sans biais.
- Mais ces stratégies optimales en vitesse n'ont pas le même coût de calcul:
Averaging : $1/\delta^2$
Nesterov : $1/\delta^{(b+1)/2}$ avec $b > 3$.
- Pour qu'il y ait convergence, il suffit
Averaging : $0 \leq c < 1$ $0 < b \leq b + c \leq 1 + a$
Nesterov : $0 \leq c < 2$ $0 < b \leq b + c < 3 - c$

Plan

Régression Logistique à Effets Aléatoires

Algorithme Gradient-Proximal

Algorithme Gradient-Proximal perturbé

Perturbation aléatoire type Monte Carlo

Régression Logistique à effets aléatoires (suite)

Mise en oeuvre des algorithmes Gradient-Proximaux stochastiques

Données simulées

Convergence de $\{\beta_n, n \geq 0\}$

Convergence de $\{F(\theta_n), n \geq 0\}$

Conclusion

Rappels

- Minimisation de la log-vraisemblance négative pénalisée

$$\min_{(\beta, \sigma) \in \mathbb{R}^p \times \mathbb{R}^+} -\ell(\theta) + g(\theta)$$

- lorsque g est convexe, sci; $-\ell$ est **non convexe**
- $\nabla \ell$ est Lipschitz et $\nabla \ell(\theta) = \int H_\theta(\mathbf{u}) \pi_\theta(\mathbf{u}) d\mathbf{u}$ avec

$$H_\theta(\mathbf{u}) = \sum_{i=1}^N (Y_i - F(x_i' \beta + \sigma z_i' \mathbf{u})) \begin{bmatrix} x_i \\ z_i' \mathbf{u} \end{bmatrix}$$

Approximation Monte Carlo

MCMC de type Gibbs, basé sur une approche “augmentation de données”

$$\nabla \ell(\theta) = \int_{\mathbb{R}^q \times \mathbb{R}^N} H_\theta(\mathbf{u}) \tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) \, d\mathbf{u} d\mathbf{w},$$

en ayant posé

$$\tilde{\pi}_\theta(\mathbf{u}, \mathbf{w}) = \left(\prod_{i=1}^N \bar{\pi}_{\text{PG}}(w_i; x'_i \beta + \sigma z'_i \mathbf{u}) \right) \pi_\theta(\mathbf{u}) d\mathbf{u}$$

où $\bar{\pi}_{\text{PG}}$ désigne la densité d'une loi Polya-Gamma.

► **Algorithme de Polson et al. (2012)** Etant donnée la valeur courante $(\mathbf{u}^t, \mathbf{w}^t)$,

(i) simuler $\mathbf{u}^{t+1} \sim \mathcal{N}_q(\mu_\theta(\mathbf{w}^t); \Gamma_\theta(\mathbf{w}^t))$

(ii) simuler $\mathbf{w}^{t+1} \sim \prod_{i=1}^N \bar{\pi}_{\text{PG}}(w_i; |\mathbf{X}_i \cdot \beta + \sigma(\mathbf{Z}\mathbf{u}^{t+1})_i|)$

$$\Gamma_\theta(\mathbf{w}) = \left(I + \sigma^2 \sum_{i=1}^N w_i \mathbf{Z}'_i \mathbf{Z}_i \right)^{-1}, \quad \mu_\theta(\mathbf{w}) = \sigma \Gamma_\theta(\mathbf{w}) \sum_{i=1}^N ((Y_i - 1/2) - w_i \mathbf{X}_i \cdot \beta) \mathbf{Z}_i.$$

Données simulées (1/2)

- $N = 500$ observations; $p = 1000$ régresseurs.
- Matrice de design \mathbf{X} construite par colonne selon un AR
 $\mathbf{X}_{\cdot, n+1} = 0.8\mathbf{X}_{\cdot, n} + \sqrt{1 - 0.8^2}\mathcal{N}_q(0, I)$.
- Coefficients de régression β_{true} : 2% de coeff non nuls.

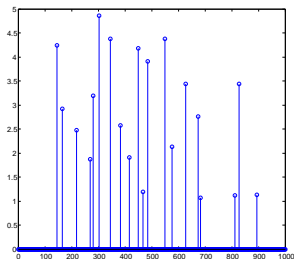
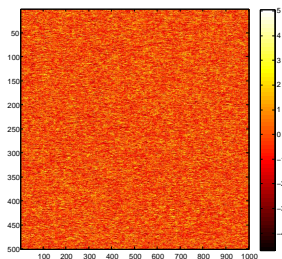


FIG.: (gauche) Matrice de design \mathbf{X} de taille 500×1000 . (droite) Coeff. de régression

Données simulées (2/2)

- $q = 5$ effets aléatoires.
- $\mathbf{U} \sim \mathcal{N}_q(0, I)$.
- Matrice \mathbf{Z} de forme simple - pour rendre qq calculs explicites dans cette illustration numérique.

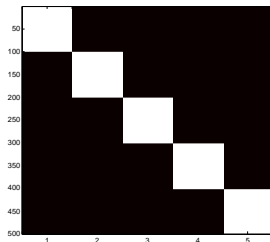


FIG.: Matrice \mathbf{Z} binaire, de taille 500×5 : 1 (en blanc) et 0 (en noir)

Convergence de β_n (1/2)

- On itère l'algorithme Gradient Proximal Stoch sur 150 itérations dans le cas $\gamma_n = 0.005$ et $m_n = 200 + n$. On relève la valeur limite β_∞ .
- On trace le support de β_∞ et celui de β_{true} pour comparaison (attention, $\beta_{\text{true}} \neq \text{argmin}(-\ell + g)$).

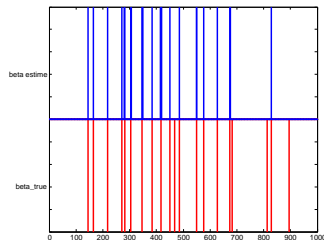


FIG.: Comparaison des supports : β_∞ (haut) et β_{true} (bas)

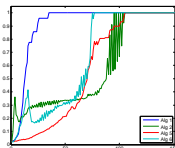
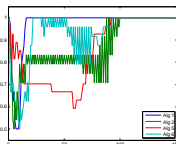
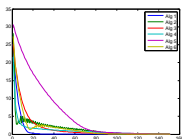
Convergence de β_n (2/2)

- On itère plusieurs algorithmes sur 150 itérations. On relève la valeur limite β_∞ .

	Alg	γ_n	m_n	autres param.
1	GP-stoch	0.005	$200 + n$	
2	GP-stoch	$0.05/\sqrt{n}$	$290 + \sqrt{n}$	
3	Averagé	0.005	$200 + n$	$a_n = \sqrt{n}$
4	Averagé	$0.05/\sqrt{n}$	$290 + \sqrt{n}$	$a_n = \sqrt{n}$
5	Nesterov	0.001	$45 + n^{3.1}/6000$	$t_n \sim n^2$
6	Nesterov	$0.005 \wedge 0.1/n$	$155 + n^{2.1}/100$	$t_n \sim n^2$

- On trace *l'erreur relative*, *la sensitivity* et *la précision*

$$\varepsilon_n = \frac{\|\beta_n - \beta_\infty\|}{\|\beta_\infty\|}, \quad \text{SEN}_n = \frac{\sum_i 1_{|\beta_{n,i}| > 0} 1_{|\beta_{\infty,i}| > 0}}{\sum_i 1_{|\beta_{\infty,i}| > 0}}, \quad \text{PRE}_n = \frac{\sum_i 1_{|\beta_{n,i}| > 0} 1_{|\beta_{\infty,i}| > 0}}{\sum_i 1_{|\beta_{n,i}| > 0}}$$



Convergence de $F(\theta_n)$

Comparaison de stratégies d'averaging. On itère plusieurs algorithmes sur 150 itérations. On trace une approximation de $n \mapsto \mathbb{E}[F(\bar{\theta}_n)] - F(\theta_\infty)$, calculée par Monte Carlo sur 50 runs indépendants. (on n'observe pas de valeurs significativement différentes pour $F(\theta_\infty)$ sur les 50 runs)

Alg	γ_n	m_n	autres param.
1	0.005	$200 + n$	$a_n = \sqrt{n}$
2	0.005	$200 + n$	$a_n = 1$
3	0.005	$200 + n$	$a_n = n^{-0.1}$
4	$0.05/\sqrt{n}$	$290 + \sqrt{n}$	$a_n = \sqrt{n}$
5	$0.05/\sqrt{n}$	$290 + \sqrt{n}$	$a_n = 1$
6	$0.05/\sqrt{n}$	$290 + \sqrt{n}$	$a_n = n^{-0.49}$

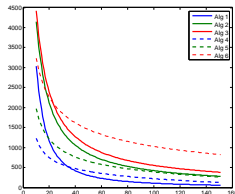
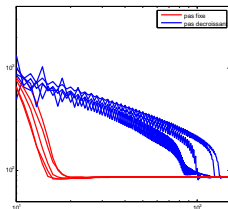


FIG.: (gauche) Trajectoires de l'algorithme Gdt Prox: à pas fixe et à pas décroissant. (droite) Estimation de $\mathbb{E}[F(\bar{\theta}_n)] - F(\theta_\star)$.

Plan

Régression Logistique à Effets Aléatoires

Algorithme Gradient-Proximal

Algorithme Gradient-Proximal perturbé

Perturbation aléatoire type Monte Carlo

Régression Logistique à effets aléatoires (suite)

Conclusion

Conclusion de la présentation

- On dispose d'outils numériques pour calculer l'estimateur du maximum de vraisemblance pénalisée dans des modèles à données latentes, dans le cas de pénalités convexes engendrant des opérateurs proximaux calculables.
- Pour les algorithmes Gradient-Proximaux stochastiques
 - la stratégie "pas fixe + mini-batch" conduit aux meilleures vitesses de convergence : $1/n$ pour des stratégies d'averaging et $1/n^2$ pour des stratégies d'accélération à la Nesterov.
 - le coût computationnel est moindre pour les algorithmes averagés : $1/\delta^2$ pour atteindre une précision de δ dans l'estimation de $\min F$.