

Exam 2016-2017

Ex. 1: (logistic regression with gross outliers)

① $f_i(x, c) = \log(1 + e^{-b_i(a_i^T x + c)})$

~~$\nabla f_i(x, c)$~~ $\nabla f_i(x, c) = \begin{pmatrix} \nabla_x f_i(x, c) \\ \nabla_c f_i(x, c) \end{pmatrix}$

$\nabla_x f_i(x, c) = - \frac{-b_i a_i e^{-b_i(a_i^T x + c)}}{(1 + e^{-b_i(a_i^T x + c)})^2} + dx = - \frac{b_i a_i}{(1 + e^{b_i(a_i^T x + c)})} + dx$

$\nabla_c f_i(x, c) = \frac{-b_i}{(1 + e^{b_i(a_i^T x + c)})}$

$$A^T \begin{pmatrix} \frac{-b_1}{1 + e^{b_1(a_1^T x + c)}} \\ \vdots \\ \frac{-b_n}{1 + e^{b_n(a_n^T x + c)}} \end{pmatrix}$$

Donc $\nabla_x f(x, c) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x, c) = - \frac{1}{n} \sum_{i=1}^n \left(\frac{b_i a_i}{(1 + e^{b_i(a_i^T x + c)})} + dx \right)$

$\nabla_c f(x, c) = - \frac{1}{n} \sum_{i=1}^n \frac{b_i}{(1 + e^{b_i(a_i^T x + c)})}$

Computing $\nabla f_i(x, c)$ is $O(d)$
 $\nabla f(x, c)$ is $O(nd)$ (for a single x)

② For $k=1, \dots, T$,
 $\left\{ \begin{array}{l} \text{sample } i \in \{1, \dots, n\} \\ x_{k+1} = x_k - \alpha \nabla f_i(x_k) \end{array} \right.$ where α = step size

③ We rewrite the iterate $x_k = \beta_k z_k$, $\beta_k \in \mathbb{R}$, $z_k \in \mathbb{R}^d$.

$x_{k+1} = x_k - \eta \left[- \frac{b_i a_i}{(1 + e^{b_i(a_i^T x_k + c)})} + dx \right]$

$= x_k (1 - d\eta) + \eta a_i l(a_i^T x_k)$ avec $l(a_i^T x_k) = \frac{b_i}{1 + e^{b_i(a_i^T x_k + c)}}$

$$\Rightarrow \beta_{k+1} z_{k+1} = (1 - d\eta) \beta_k z_k + \eta a_i l(a_i^T \beta_k z_k) \quad O(1)$$

$$= \underbrace{(1 - d\eta) \beta_k}_{\beta_{k+1}} \left[z_k - \frac{\eta a_i l(a_i^T \beta_k z_k)}{(1 - d\eta) \beta_k} \right] \quad O(1)$$

$$\beta_{k+1} = (1 - d\eta) \beta_k \quad \{O(1)\}$$

$$z_{k+1} = z_k - \frac{\eta a_i l(a_i^T \beta_k z_k)}{(1 - d\eta) \beta_k} \quad \{O(1)\}$$

④ c'est le vecteur c of the biases is of length \mathbb{R}^n . But we want to shrink the c_i of the outliers?

We use the l_1 norm because the l_2 is more robust to outliers. The higher γ , the more c_i are set to zero.

⑤ à vérifier

⑥ On peut réécrire $f(x, c)$ par: $\log(1 + \exp(-b_i(a_i^T x + c_i))) + \frac{1}{2} \|x\|_2^2$
Donc $\nabla_x f(x, c)$ ne change pas et est toujours égal à:

$$\frac{-b_i c_i}{1 + e^{b_i(a_i^T x + c_i)}} + \frac{1}{2} x$$

et: $\nabla_c f(x, c) = \frac{-b_i c_i}{1 + e^{b_i(a_i^T x + c_i)}} + \frac{1}{2} x$

Donc $\nabla f(x, c) \in \mathbb{R}^{d+n}$ peut s'écrire comme:

$$x_i \left\{ \frac{b_i}{1 + e^{b_i(a_i^T x + c_i)}} \times \begin{pmatrix} a_i - a_i + \frac{1}{2} x \\ -c_i \end{pmatrix} \right\} := \beta_i \in \mathbb{R}^d$$

et $\nabla f(x, c) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x, c) = \frac{1}{n} \sum_{i=1}^n \beta_i$

⑦ Algo:

$$\begin{cases} x_{t+1} = x_t - \eta \left(\frac{1}{n} \sum_{i=1}^n \frac{-b_i a_i}{1 + e^{b_i(a_i^T x_t + c_t)}} + \frac{1}{2} x_t \right) \\ c^{t+1} = \text{proj}_{\mathcal{C}} \left(c_t + \frac{\eta \sum_{i=1}^n b_i c_i}{1 + e^{b_i(a_i^T x_t + c_t)}} \right) \\ \text{Avec } \text{proj}_{\mathcal{C}} \|c\|_1 = (S_{\eta \gamma}(c_1), \dots, S_{\eta \gamma}(c_n)) \end{cases}$$

Exo.2: Elements of proof of the SAGA algorithm

$$\begin{aligned} \textcircled{1} T^k &= \frac{1}{n} \sum_{i=1}^n f_i(\Phi_i^k) - f(x^*) - \frac{1}{n} \sum_{i=1}^n \langle \nabla f_i(x^*), \Phi_i^k - x^* \rangle + c \|x^k - x^*\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left(\underbrace{f_i(\Phi_i^k) - \langle \nabla f_i(x^*), \Phi_i^k - x^* \rangle}_{\text{zero here convex}} \right) + c \|x^k - x^*\|^2 \end{aligned}$$

$$\text{Donc } T^k \geq c \|x^k - x^*\|^2$$

$\textcircled{2}$ According to the previous question, $T^k \geq c \|x^k - x^*\|^2$
Therefore, $c \mathbb{E}[\|x^k - x^*\|^2] \leq \mathbb{E}[T^k]$

$$\begin{aligned} &= \mathbb{E}[\mathbb{E}(T^k | \mathcal{F}_{k-1})] \\ &\leq \mathbb{E}\left[\left(1 - \frac{1}{K}\right) T^{k-1}\right] \quad \text{d'après la question 1} \\ &\leq \left(1 - \frac{1}{K}\right) \mathbb{E}(T^{k-2}) \\ &= \left(1 - \frac{1}{K}\right) \mathbb{E}[\mathbb{E}(T^{k-1} | \mathcal{F}_{k-2})] \\ &\leq \left(1 - \frac{1}{K}\right)^2 \mathbb{E}(T^{k-4}) \end{aligned}$$

$$\leq \left(1 - \frac{1}{K}\right)^k \mathbb{E}(T_0)$$

lasto est déterminé, donc: $\mathbb{E}(T_0) = T_0$ et:

$$\boxed{c \mathbb{E}[\|x^k - x^*\|^2] \leq \left(1 - \frac{1}{K}\right)^k T_0}$$

~~Donc $\mathbb{E}[\|x^k - x^*\|^2] \leq \frac{T_0}{c} \left(1 - \frac{1}{K}\right)^k$~~

Let's now prove that (1) holds:

$$T_0 = \frac{1}{n} \sum_{i=1}^n f_i(\Phi_i^0) - f(x^*) - \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^*)^T (\Phi_i^0 - x^*) + c \|x_0 - x^*\|^2$$

$$= f(x_0) - f(x^*) - \nabla f(x^*)^T (x_0 - x^*) + c \|x_0 - x^*\|^2$$

$$\text{Donc } \mathbb{E}[\|x^k - x^*\|^2] \leq \left(\frac{1}{c} (f(x_0) - \nabla f(x^*)^T (x_0 - x^*) - f(x^*)) + \|x_0 - x^*\|^2 \right) \left(1 - \frac{1}{K}\right)^k$$

D'après l'énoncé, ~~$\frac{1}{c} (f(x_0) - \nabla f(x^*)^T (x_0 - x^*) - f(x^*)) \leq \frac{\mu}{2(\eta\mu + 1)}$~~

$$\frac{1}{c} = \frac{\mu}{2(\eta\mu + 1)} \Rightarrow 1 - \frac{1}{K} = 1 - \frac{\mu}{2(\eta\mu + 1)}$$

$$\begin{aligned} \text{et } 1 - \frac{1}{c} &= 2\sigma n(1 - \delta_\mu) = \frac{R}{\eta} \frac{n}{(\eta\mu + 1)} \left(1 - \frac{\mu}{2(\eta\mu + 1)}\right) \\ &\leq \frac{n}{\eta\mu + 1} \leq 1 \end{aligned}$$

En concert, $f(x^0) - \langle \nabla f(x^*), x^0 - x^* \rangle \geq 0$
 (car $f(x^*) = 0$)

Donc: $\frac{1}{c} (\quad) \leq \frac{1}{1+L} (\quad)$

→ donc on a prouvé (1).

(3) $\mathbb{E}[T^{k+1} | \mathcal{F}_k] = \mathbb{E}_{\mathcal{I}_k}[T^{k+1}]$
 $= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{I}_k}[f_i(\phi_i^{k+1})] - f(x^*) - \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{I}_k}[\nabla f_i(x^*)^T (\phi_i^k - x^*)]}_B + c \mathbb{E}_{\mathcal{I}_k}[\|x^k - x^*\|]$

$\mathbb{E}_{\mathcal{I}_k}[f_i(\phi_i^{k+1})] = \frac{1}{n} f_i(x^k) + (1 - \frac{1}{n}) f_i(\phi_i^k)$
 (par définition)

Donc: $\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{I}_k}[f_i(\phi_i^{k+1})] = \frac{1}{n} \left(\sum_{i=1}^n \left(\frac{1}{n} f_i(x^k) + (1 - \frac{1}{n}) f_i(\phi_i^k) \right) \right)$
 $= \frac{1}{n} f(x^k) + (1 - \frac{1}{n}) \frac{1}{n} \sum_{i=1}^n f_i(\phi_i^k)$

De même, $\mathbb{E}_{\mathcal{I}_k}[\nabla f_i(x^*)^T (\phi_i^{k+1} - x^*)]$
 $= \frac{1}{n} \nabla f_i(x^*)^T (x^k - x^*) + (1 - \frac{1}{n}) \nabla f_i(x^*)^T (\phi_i^k - x^*)$

Donc $B = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{I}_k}[\nabla f_i(x^*)^T (\phi_i^{k+1} - x^*)]$
 $= \frac{1}{n} \sum_{i=1}^n \frac{1}{n} \nabla f_i(x^*)^T (x^k - x^*) + (1 - \frac{1}{n}) \nabla f_i(x^*)^T (\phi_i^k - x^*)$
 $= \frac{1}{n} \nabla f(x^*)^T (x^k - x^*) + (1 - \frac{1}{n}) \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^*)^T (\phi_i^k - x^*)$

Donc $-\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{I}_k}[\nabla f_i(x^*)^T (\phi_i^{k+1} - x^*)]$
 $= -\frac{1}{n} \nabla f(x^*)^T (x^k - x^*) - \frac{1}{n} (1 - \frac{1}{n}) \sum_{i=1}^n \nabla f_i(x^*)^T (\phi_i^k - x^*)$

En reliant les morceaux, on a la double inégalité.

④ D'après l'énoncé, $\phi_{I_k}^{k+1} = x^k$.

$$\begin{aligned} \mathbb{E}_{I_k} [\nabla f_{I_k}(\phi_{I_k}^{k+1})] &= \mathbb{E}_{I_k} [\nabla f_{I_k}(x^k)] \\ &= \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) = \nabla f(x^k). \end{aligned}$$

Pour prouver l'égalité, il reste à montrer que :

$$\mathbb{E} [-\nabla f_{I_k}(x^k) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k)] = 0.$$

$$\text{On a : } \mathbb{E} \left[-\nabla f_{I_k}(x^k) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) \right] = -\frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) = 0$$

⑤ $\|x^{k+1} - x^*\|^2 = \|x^{k+1} - x^k + x^k - x^*\|^2$

$$= \|x^k - x^*\|^2 + 2(x^{k+1} - x^k)^T (x^k - x^*) + \|x^{k+1} - x^k\|^2$$

mais $x^{k+1} - x^k = -\gamma \left(\nabla f_{I_k}(x^k) - \nabla f_{I_k}(\phi_{I_k}^k) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\phi_i^k) \right)$

$$\begin{aligned} \mathbb{E}_{I_k} [\|x^{k+1} - x^*\|^2] &= \underbrace{\mathbb{E}_{I_k} [\|x^k - x^*\|^2]}_{\text{déterministe}} + 2(x^k - x^*)^T \mathbb{E}_{I_k} (x^{k+1} - x^k) \\ &\quad + \mathbb{E}_{I_k} (\|x^{k+1} - x^k\|^2). \end{aligned}$$

$$= \|x^k - x^*\|^2 - 2\gamma(x^k - x^*)^T \nabla f(x^k) + \gamma^2 \mathbb{E}_{I_k} \left(\left\| \nabla f_{I_k}(x^k) - \nabla f_{I_k}(\phi_{I_k}^k) + \frac{1}{n} \sum_{i=1}^n \nabla f_i(\phi_i^k) \right\|^2 \right)$$

Comme f est convexe, et on minimise sur \mathbb{R}^d alors $\nabla f(x^*) = 0$
Donc le résultat.

⑥

⑦ Prouver que g est $L-\mu$ smooth.

$$g(y) - g(x) = f(y) - f(x) - \frac{\mu}{2} \|y\|^2 + \frac{\mu}{2} \|x\|^2$$

$$\leq \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2 - \frac{\mu}{2} \|y\|^2 + \frac{\mu}{2} \|x\|^2$$

g is L smooth

Mais $\nabla f(x)^T = (\nabla g(x) + \mu x)^T$

Donc: $g(y) - g(x) \leq \nabla g(x)^T (y-x) + \mu x^T (y-x) - \frac{\mu}{2} \|y\|^2 + \frac{\mu}{2} \|x\|^2$

Mais $\|y\|^2 = \|y-x+x\|^2 = \|y-x\|^2 + \|x\|^2 + 2(y-x)^T x$

$\Rightarrow -\frac{\mu}{2} \|y\|^2 = -\frac{\mu}{2} \|y-x\|^2 - \frac{\mu}{2} \|x\|^2 - \mu x^T (y-x)$

Donc $g(y) - g(x) \leq \nabla g(x)^T (y-x) + \frac{L}{2} \|y-x\|^2 - \frac{\mu}{2} \|y-x\|^2$

$$= \nabla g(x)^T (y-x) + \frac{L-\mu}{2} \|y-x\|^2$$

Donc g est $(L-\mu)$ smooth.

Donc: $\Rightarrow \nabla g(y) - \nabla g(x) \leq \frac{1}{L-\mu} \|\nabla g(y) - \nabla g(x)\|$