

Starting 10/5

Note on homeworks:

Please only submit standard file formats

.pdf, .zip, images (jpeg, png, etc).

HW1 posted grades (shortly)

Questions about grading details, see corresponding TA:

$\begin{cases} Q1 \leftarrow \text{Jinke} \\ Q2 \leftarrow \text{Thabo} \\ Q3 \leftarrow \text{Ning} \\ Q4 \leftarrow \text{Salma} \end{cases} \quad (\sim 16/20)$

### 3. Projected/Proximal Gradient Descent

Two examples currently motivating us ...

LASSO  $\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 + \gamma \|x\|_1$

Grading Program  $\min_{x \in \mathbb{R}^d} f(x) + \delta_S(x)$

Lets consider the more general model

$\min_{x \in \mathbb{R}^d} f(x) + h(x)$

$\uparrow$  smooth nonconvex(?)       $\uparrow$  convex, "structured" (i.e. we can compute  $\text{prox}_h(\cdot)$ )

Lemma (Sum Rule)  $\partial(f+h)(x) = \partial f(x) + \partial h(x)$   
if  $f, h$  convex.

Proof. " $\supseteq$ " Let  $g_1 \in \partial f(x), g_2 \in \partial h(x)$

$$\begin{aligned} f(y) &\geq f(x) + g_1^T(y-x) \quad \forall y \\ + \quad h(y) &\geq h(x) + g_2^T(y-x) \quad \forall y \end{aligned}$$

$$\hline (f+h)(y) \geq (f+h)(x) + (g_1 + g_2)^T(y-x) \quad \forall y. \quad \checkmark$$

" $\subseteq$ " Harder to prove, Need Separating Hyperplane Thm.  
("Intro to Convexity")  $\square$

~~Lemma~~ For particular problem, we consider  
 $\nabla f(x) + \partial h(x).$

(Proximal First-Order Condition).

$\angle_{\text{prox}}$  is our angle of attack.

only need  $C'$

Lemma For any  $f$  with  $L$ -Lipschitz gradient and  $h$  convex,  
if  $x^*$  is a local min of  $f+h$ , then  
 $0 \in \nabla f(x^*) + \partial h(x^*)$ .

Proof. Suppose  $-\nabla f(x^*) \notin \partial h(x^*)$ .

$\Rightarrow$  Some  $y \in \mathbb{R}^d$  s.t.  $h(y) < \underline{h(x^*) + -\nabla f(x^*)^T(y-x^*)}$

Consider  $z = x^* + \lambda(y-x^*)$  for  $0 \leq \lambda \leq 1$ .

$$\begin{aligned} h(z) &\leq (1-\lambda)h(x^*) + \lambda h(y) \\ &= h(x^*) + \lambda(h(y) - h(x^*)) \\ &< h(x^*) - \lambda \nabla f(x^*)^T(y-x^*) \end{aligned}$$

$$\text{As } \lambda \rightarrow 0, z \rightarrow x^*, \quad \nabla f(x^*)^T(y-x^*) \leq \frac{h(z) - h(x^*)}{\lambda}$$

$$\Rightarrow h(z) < h(x^*) - \lambda \cdot \frac{h(z) - h(x^*)}{\lambda}$$

$$\Rightarrow (f+h)(z) < (f+h)(x^*)$$

$\Rightarrow z$  is better than  $x^*$

$\Rightarrow x^*$  is not a local min.  $\square$

Lets develop an iterative algorithm to min  $f+h$ :

At step  $k$ , we can linearize  $f$

$$(f+h)(x) \approx \underbrace{f(x_k) + \nabla f(x_k)^T(x-x_k)}_{\substack{\text{1st-order model} \\ \text{of } f \text{ at } x_k}} + h(x)$$

Lets iterate

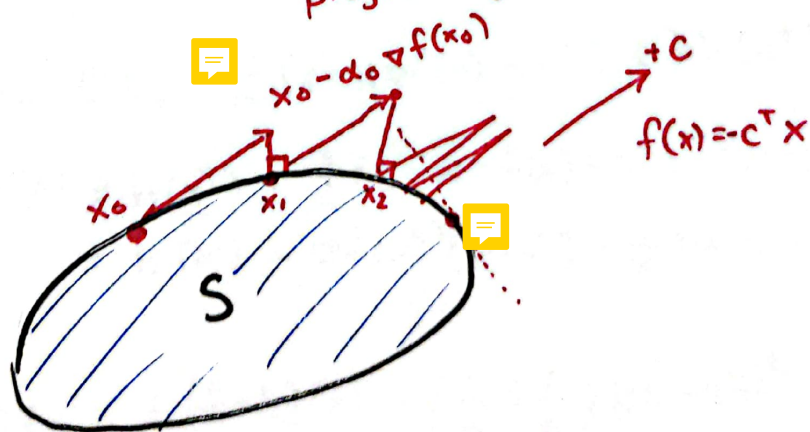
$$\begin{aligned}
 x_{k+1} &= \text{prox}_{\alpha_k} (f(x_k) + \nabla f(x_k)^T (x - x_k) + h(x)) (x_k) \\
 &= \arg \min_{x \in \mathbb{R}^d} \left\{ f(x_k) + \nabla f(x_k)^T (x - x_k) + h(x) + \frac{1}{2\alpha_k} \|x - x_k\|_2^2 \right\} \\
 &= \arg \min \left\{ \underbrace{f(x_k) - \frac{1}{2\alpha_k} \|\nabla f(x_k)\|^2}_{\text{constant}} + h(x) + \frac{1}{2\alpha_k} \|x - (x_k - \alpha_k \nabla f(x_k))\|^2 \right\} \\
 &= \arg \min \left\{ h(x) + \frac{1}{2\alpha_k} \|x - (x_k - \alpha_k \nabla f(x_k))\|^2 \right\} \\
 &= \text{prox}_{\alpha_k h} (x_k - \alpha_k \nabla f(x_k)).
 \end{aligned}$$

This alternates grad descent on  $f$ , prox on  $h$ , repeat....

Called "projected/proximal gradient descent"  
"ISTA"

"Forward-Backward Method"

projected gradient descent.





Define  $G_\alpha(x) = \frac{1}{\alpha}(x - \underbrace{\text{prox}_{\alpha h}(x - \alpha \nabla f(x))}_{x^*})$   
 as the gradient mapping.

Check "gradient-like". HW 2, Q2.

$$\frac{1}{\alpha}(x - \alpha \nabla f(x) - \text{prox}_{\alpha h}(x - \alpha \nabla f(x))) \in \partial h(x^*)$$

$$\Leftrightarrow G_\alpha(x) \in \nabla f(x) + \partial h(x^*)$$

Small  $G_\alpha(x) \Rightarrow x, x^*$  are essentially the same.

Lemma (Descent) For any  $x$ , let  $x^* = \text{prox}_{\alpha h}(x - \alpha \nabla f(x))$ .

$$(f+h)(x^*) \leq (f+h)(x) - \left(\alpha - \frac{L\alpha^2}{2}\right) \|G_\alpha(x)\|^2$$

whenever  $\nabla f$  is  $L$ -Lipschitz.

Proof. By our Taylor Approximation Thms,

$$f(x^*) \leq f(x) + \nabla f(x)^T(x^* - x) + \frac{L}{2} \|x^* - x\|^2. \quad (1)$$

By HW 2, Q2,  $\frac{1}{\alpha}(x - \alpha \nabla f(x) - x^*) \in \partial h(x^*)$

$$\Rightarrow h(x) \geq h(x^*) + \frac{1}{\alpha}(x - \alpha \nabla f(x) - x^*)^T(x - x^*)$$

$$= h(x^*) - \nabla f(x)^T(x - x^*) + \frac{1}{\alpha} \|x - x^*\|^2. \quad (2)$$

(1) - (2)

$$(f+h)(x^*) \leq (f+h)(x) - \left(\frac{1}{\alpha} - \frac{L}{2}\right) \|x - x^*\|^2$$

$$= (f+h)(x) - \left(\alpha - \frac{L\alpha^2}{2}\right) \|G_\alpha(x)\|^2.$$

Picking  $\alpha = \frac{1}{L}$  gets descent

$$(f+h)(x^+) \leq (f+h)(x) - \frac{1}{2L} \|G_{1/2}(x)\|^2$$

Linesearching (exact, backtracking) work exactly the same.

(Beck Ch 10 repeats these for us).

Theorem For any  $f$  with  $L$ -Lipschitz gradient and convex  $h$ , selecting  $\alpha_k = \frac{1}{L}$ , the proximal gradient method

has

$$\frac{1}{T} \sum_{k=0}^{T-1} \|G_k(x_k)\|^2 \leq \frac{2L(f+h)(x_0) - \min f+h}{T}.$$

Proof. Our descent lemma at each iteration gives

$$(f+h)(x_{k+1}) \leq (f+h)(x_k) - \frac{1}{2L} \|G_{1/2}(x_k)\|^2.$$