

Subgradiates y subdiferenciales

Dada una función convexa $f: \mathbb{R}^n \rightarrow \mathbb{R}$, decimos que un vector $d \in \mathbb{R}^n$ es un subgradiente de f en un punto x si

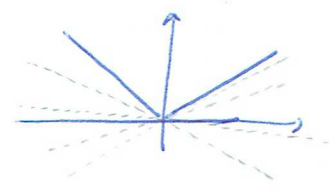
$$f(z) \geq f(x) + (z-x)^T d \quad \forall z \in \mathbb{R}^n$$

Al conjunto de todos los subgradiates de f en x se le llama subdiferencial de f en x , y se denota $df(x)$.

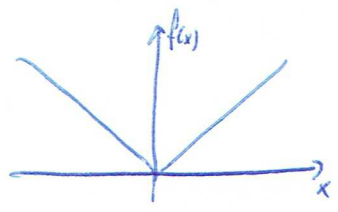
Ejemplos

$f(x) = |x|$ en $x = 0$

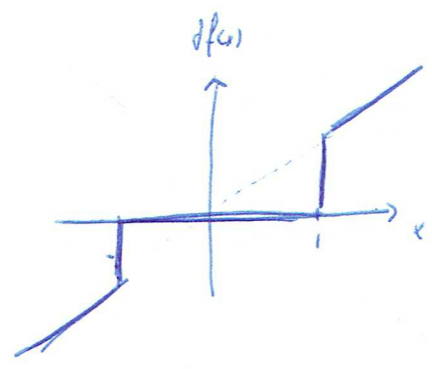
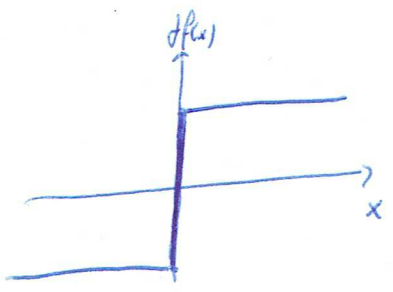
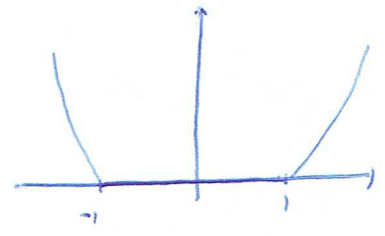
$df(0) = [-1, 1]$



$f(x) = |x|$



$f(x) = \max\{0, \frac{1}{2}(x^2-1)\}$



Obs

Si f es diferenciable, el único subgradiente es el gradiente.

Prop

Sea $f: \mathbb{R}^n \rightarrow \mathbb{R}$ convexa, $x \in \mathbb{R}^n$. Entonces $df(x)$ es un conjunto no vacío, convexo y compacto. Además, $\forall w \in \mathbb{R}^n$ se tiene

$$\frac{df}{dw} = \max_{g \in df(x)} g^T w \quad \left(\text{con } \frac{df}{dw} = \lim_{\alpha > 0} \frac{f(x+\alpha d) - f(x)}{\alpha} \right)$$

Dem (de una parte)

Por la definición de subgradiente, tenemos $\frac{f(x+\alpha d) - f(x)}{\alpha} \geq g^T d \quad \forall d \in \mathbb{R}^n, \alpha > 0$

Como el lado izquierdo converge monótonamente (con $\alpha > 0$) a $\frac{df}{dw}$, entonces

$$g \in df(x) \iff \frac{df}{dw} \geq g^T d \quad \forall d \in \mathbb{R}^n$$

Por lo tanto $df(x)$ se puede escribir como la intersección de los conjuntos

$$\left\{ g \mid \frac{df}{dw} \geq g^T d \right\} \quad \text{haciendo variar } d \text{ a todos los vectores no nulos de } \mathbb{R}^n.$$

Estos conjuntos son semi-espacios cerrados, y por lo tanto $df(x)$ es cerrado y convexo.

Además si fuera no acotado, podríamos hacer $g^T d$ no acotado, por el α d , lo que contradice $\frac{df}{dw} \geq g^T d$

Prop (Condición de Optimalidad)

Un punto $x \in \mathbb{R}^n$ minimiza una función convexa $f: \mathbb{R}^n \rightarrow \mathbb{R}$ sobre un conjunto convexo X si existe un subgradiente $g \in df(x)$ tal que $g^T(z-x) \geq 0 \quad \forall z \in X$

Obs:

- Esta condición generaliza la condición $\nabla f(x)^T(z-x) \geq 0 \quad \forall z \in X$
- Cuando $X = \mathbb{R}^n$, la condición necesaria y suficiente para $0 \in df(x)$

Interpretación geométrica

Dado un conjunto X y $x \in X$, llamamos cono normal a X en x al conjunto

$$N_x(x) = \{ \varphi \mid \varphi^T(z-x) \leq 0, \forall z \in X \}$$

Subgradiante

Fijemos $x \in \mathbb{R}^n$, y llamemos $c = f(x)$

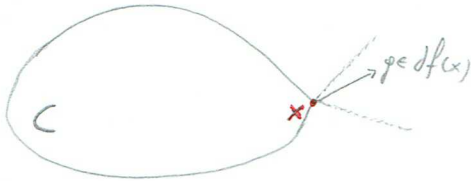
Consideremos el conjunto de nivel $C = \{ z \in \mathbb{R}^n \mid f(z) \leq c \}$

$$\text{Si } \varphi \in \partial f(x) \Rightarrow f(z) \geq f(x) + \varphi^T(z-x) \quad \forall z \in \mathbb{R}^n$$

$$\text{En particular, si } z \in C \Rightarrow f(z) \leq f(x) \Rightarrow \varphi^T(z-x) \leq 0$$

Esto quiere decir que $\varphi \in N_C(x)$

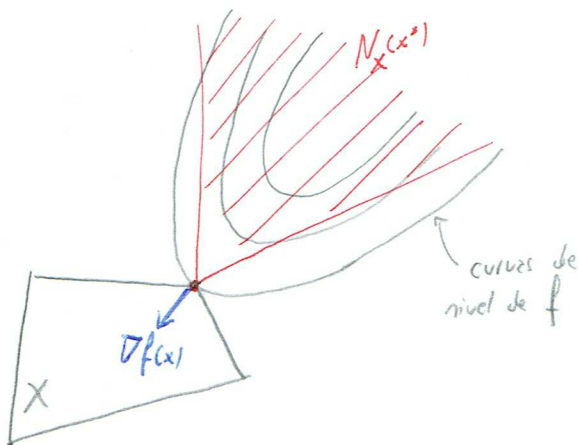
(Esto generaliza el concepto que el subgradiante es normal a la curva de nivel. Cuando f es diferenciable, el cono normal a C es una semi-recta)



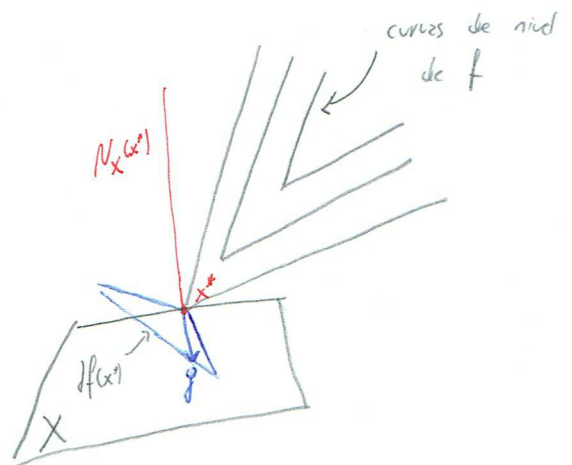
Condición de optimalidad

La condición de optimalidad ($\exists \varphi \in \partial f(x) \mid \varphi^T(z-x) \geq 0 \quad \forall z \in X$)

se traduce en que $\exists \varphi \in \partial f(x) \mid -\varphi \in N_x(x)$



f diferenciable \Rightarrow la condición es $-\nabla f(x^*) \in N_x(x^*)$



f no diferenciable

Métodos de subgradiente

Vemos métodos de la forma:

$$x^{k+1} = P_X(x^k - \alpha^k g^k)$$

donde g^k es cualquier subgradiente $g^k \in \partial f(x^k)$

$$\alpha^k \geq 0$$

Cuando $X = \mathbb{R}^n$ (sin restricciones)

$$x^{k+1} = x^k - \alpha^k g^k$$

Ejemplo

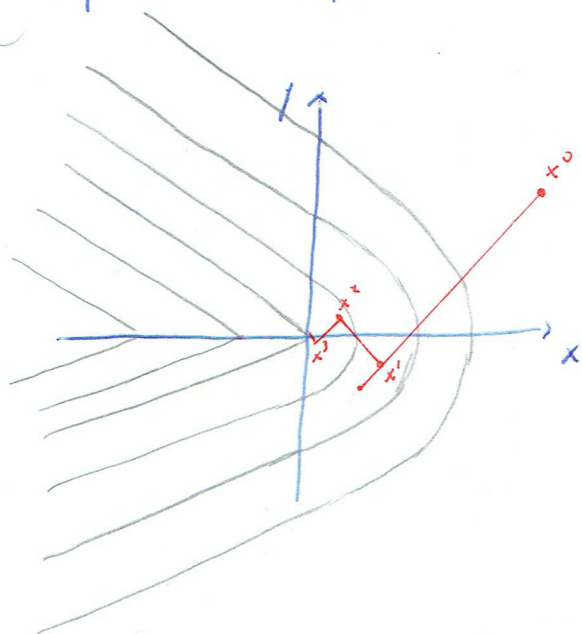
$$f(x,y) = \begin{cases} 5\sqrt{9x^2+16y^2} & \text{si } x > |y| \\ 9|x+16| & \text{si } x \leq |y| \end{cases}$$

Empezando en cualquier punto de $\{(x,y) \mid x > |y| > (\frac{9}{16})^2 x\}$

y usando el método de máximo descenso, con α^k determinado como

$$\alpha^k = \operatorname{arg\,min}_{\alpha > 0} f(x^k - \alpha g^k) \quad (\text{se puede encontrar este } \alpha^k \text{ analíticamente})$$

el método sigue una poligonal con ángulos rectos sucesivos, se converge al $(0,0)$ que no es un punto estacionario. (y de hecho $\lim_{x \rightarrow -\infty} f(x,0) = -\infty$)



A pesar de que en todos los puntos x^k la f es diferenciable, el algoritmo falla

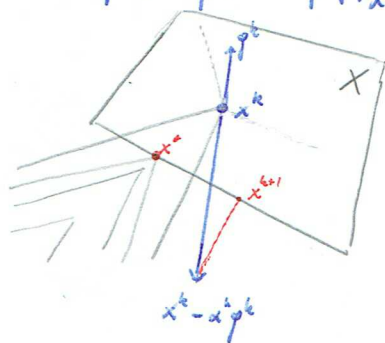
Obs:

• Al proyectar, no nos alejamos del óptimo:

De la no expansividad de la proyección: $\|P_X(x^k - \alpha^k p^k) - x^*\| \leq \|x^k - \alpha^k p^k - x^*\| \quad \forall x \in X$
en particular por x^*

• No tenemos asegurado que en cada iteración bajemos el costo.

Puede pasar que $f(P_X(x^k - \alpha^k p^k)) > f(x^k) \quad \forall \alpha > 0$



Sin embargo, si el peso es suficientemente chico, la distancia al óptimo se reduce:

Prop

Sea $\{x^k\}$ generado por el método de subgradiente. Entonces $\forall y \in X, k \geq 0$:

(a) $\|x^{k+1} - y\|^2 \leq \|x^k - y\|^2 - 2\alpha^k (f(x^k) - f(y)) + \alpha^{k2} \|p^k\|^2$

(b) si $f(y) < f(x^k)$, entonces $\|x^{k+1} - y\| < \|x^k - y\|$

por todo peso $\alpha^k / \quad 0 < \alpha^k < \frac{2(f(x^k) - f(y))}{\|p^k\|^2}$

Dem

$$\|x^{k+1} - y\|^2 = \|P_X(x^k - \alpha^k p^k) - y\|^2 \leq \|x^k - \alpha^k p^k - y\|^2 = \|x^k - y\|^2 - 2\alpha^k p^{kT}(x^k - y) + \alpha^{k2} \|p^k\|^2$$

$$\leq \|x^k - y\|^2 - 2\alpha^k (f(x^k) - f(y)) + \alpha^{k2} \|p^k\|^2$$

(en la última desigualdad se usa la definición de subgradiente)

(b) es inmediato de (a)

Esto supone tomar un paso $\alpha^k = \frac{f(x^k) - f(x^*)}{\|p^k\|^2}$ (punto medio del intervalo de la prop. anterior)

Sin embargo, muy raramente conocemos el valor óptimo.

Veamos algunos resultados por pasos "prácticos."

Paso constante:

Tomemos $\alpha^k = \alpha \ \forall k$, sea $f_\infty = \liminf_{k \rightarrow \infty} f(x^k)$

Si $\|p^k\| \leq c \ \forall k$, entonces $f_\infty \leq f(x^*) + \frac{\alpha c^2}{2}$

Dem

Si no fuera cierto, entonces $\exists \epsilon > 0 / f_\infty \geq f(x^*) + \frac{\alpha c^2}{2} + 2\epsilon$

Sea \bar{k} suficientemente grande tal que $f(x^k) \geq f_\infty - \epsilon \ \forall k \geq \bar{k}$

Sumando: $f(x^k) - f(x^*) \geq \frac{\alpha c^2}{2} + \epsilon$

Usando la Prop anterior:

$$\begin{aligned} \|x^{k+1} - x^*\|^2 &\leq \|x^k - x^*\|^2 - 2\alpha (f(x^k) - f(x^*)) + \alpha^2 c^2 \\ &\leq \|x^k - x^*\|^2 - 2\alpha \left(\frac{\alpha c^2}{2} + \epsilon\right) + \alpha^2 c^2 = \|x^k - x^*\|^2 - 2\alpha \epsilon \end{aligned}$$

Entonces, $\|x^{k+1} - x^*\|^2 \leq \|x^k - x^*\|^2 - 2\alpha \epsilon$

$$\leq \|x^{k-1} - x^*\|^2 - 4\alpha \epsilon$$

⋮

$$\leq \|x^{\bar{k}} - x^*\|^2 - 2(k+1 - \bar{k})\alpha \epsilon$$

lo cual no puede ser cierto para k grande.

Es decir, la distancia al óptimo no se puede reducir una cantidad constante ($2\alpha \epsilon$) indefinidamente. Como la distancia al óptimo está relacionada con la diferencia de valores funcionales (por la Prop anterior), llegamos a lo deseado.

Peso decreciente

(27)

Supongamos que $\|g^k\| \leq c \alpha^k$.

Entonces si $\lim_{k \rightarrow \infty} \alpha^k = 0$ y $\sum \alpha^k = \infty \Rightarrow f_\infty = f^*$

Si además $\sum (\alpha^k)^2 < \infty \Rightarrow x^k$ converge al óptimo (o a un óptimo)

La demostración es igual que el caso anterior.