

# Proximal Methods

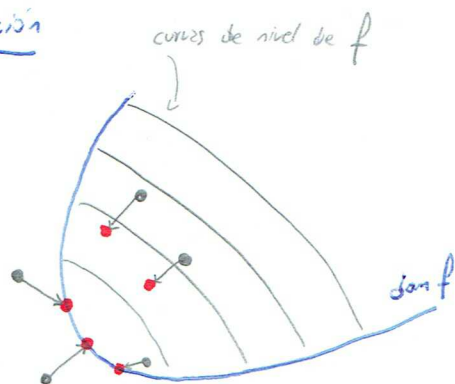
Sea  $f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$  convexa y propia ( $f \neq +\infty$ )

Definimos el operador proximal  $\text{prox}_{\lambda f}: \mathbb{R}^n \rightarrow \mathbb{R}^n$  como

$$\text{prox}_{\lambda f}(z) = \underset{x}{\text{argmin}} \left( f(x) + \frac{1}{2\lambda} \|x - z\|_2^2 \right)$$

La función  $f(x) + \frac{1}{2\lambda} \|x - z\|_2^2$  es estrictamente convexa, por lo que el mínimo es único, y la función  $\text{prox}_{\lambda f}(z)$  está bien definida.

## Interpretación



El parámetro  $\lambda$  repite qué tanto se mueve hacia el mínimo de  $f$ , y cuánto se queda cerca de  $z$ .

Obs: Si  $f$  es la función indicativa de un conjunto  $C$  cerrado y convexo:

$$f(x) = \begin{cases} 0 & x \in C \\ +\infty & x \notin C \end{cases}, \text{ entonces el operador proximal coincide con la proyección a } C.$$

## Propiedad (Separabilidad de la suma)

$$\text{Si } f(x, y) = \psi(x) + \varphi(y) \Rightarrow \text{prox}_f(w, w) = (\text{prox}_{\psi}(w), \text{prox}_{\varphi}(w))$$

$$\text{En particular si } f(x) = \sum_{i=1}^2 f_i(x_i) \Rightarrow (\text{prox}_f(w))_i = \text{prox}_{f_i}(w_i)$$

## Prop (del punto fijo)

$$x^* \text{ minimize } f \text{ sii } x^* = \text{prox}_f(x^*) \quad (\text{consideremos } \lambda=1 \text{ sin pérdida de generalidad})$$

Dem

$$\Rightarrow x^* \text{ minimiza } f \Rightarrow f(x) \geq f(x^*) \quad \forall x \Rightarrow f(x) + \frac{1}{2} \|x - x^*\|_2^2 \geq f(x^*) + \frac{1}{2} \|x^* - x^*\|_2^2 \quad \forall x$$

$$\Rightarrow x^* \text{ minimiza } f(x) + \frac{1}{2} \|x - x^*\|_2^2 \Rightarrow x^* = \text{prox}_f(x^*)$$

$$\Leftrightarrow y \text{ minimiza } f(x) + \frac{1}{2} \|x - z\|_2^2 \quad (\text{es decir } y = \text{prox}_f(z)) \quad \text{si } 0 \in \partial f(y) + (y - z)$$

$$\text{tomando } y = z = x^* \text{ tenemos que } 0 \in \partial f(x^*) \Rightarrow x^* \text{ minimiza } f$$

Entonces buscar el mínimo de  $f$  es equivalente a buscar puntos fijos del operador proximal.

Si  $\text{prox}_f$  fuera una contracción  $\Rightarrow x^{k+1} = \text{prox}_f(x^k)$  converge a un punto fijo.

Sin embargo  $\text{prox}_f$  no es necesariamente una contracción.

Veremos más adelante que de todas maneras tenemos resultados de convergencia.

### Ejemplo

Tomemos  $f(x) = |x|$  y calculemos  $\text{prox}_f(w)$

La condición de optimalidad es  $0 \in \partial f(x) + x - w \Leftrightarrow w \in x + \partial f(x)$

Dado  $w$ , para calcular  $x = \text{prox}_f(w)$  tenemos que encontrar el  $x$  que verifique la cond. de optimalidad  $w \in x + \partial f(x)$

$$\text{Entonces, si } w > 1 \Rightarrow x = w - 1$$

$$\text{si } w < -1 \Rightarrow x = w + 1$$

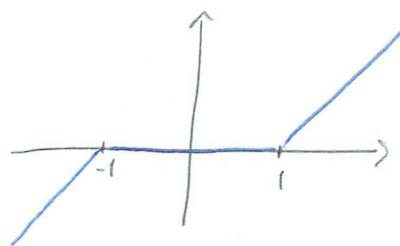
$$\text{si } |w| < 1 \Rightarrow x = 0$$

$$\text{o sea } \text{prox}_f(w) = \begin{cases} w-1 & \text{si } w > 1 \\ 0 & \text{si } |w| < 1 \\ w+1 & \text{si } w < -1 \end{cases}$$

Por la prop. de separabilidad, si  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$f(x) = \|x\|_1$ , entonces

$$(\text{prox}_f(w))_i = \begin{cases} w_i - 1 & \text{si } w_i > 1 \\ 0 & \text{si } |w_i| < 1 \\ w_i + 1 & \text{si } w_i < -1 \end{cases}$$



## Interpretación: Gradient Flow

Consideremos la ecuación diferencial  $\dot{x}(t) = -\nabla f(x(t))$  (flujo de gradiente)

Los puntos de equilibrio son ceros de  $\nabla f$ , que es donde se minimiza  $f$ .

Discretizando la ec. dif. (método Forward-Euler):

$$\frac{x^{k+1} - x^k}{h} = -\nabla f(x^k) \Rightarrow x^{k+1} = x^k - h \nabla f(x^k)$$

que es el método de descenso por gradiente con paso  $h$ .

Ahora, si usamos Backward-Euler:

$$\frac{x^{k+1} - x^k}{h} = -\nabla f(x^{k+1})$$

que tiene mejores propiedades de convergencia en ec. dif.

resulta  $x^{k+1} + h \nabla f(x^{k+1}) = x^k$  que es exactamente la condición de optimalidad por  $x^{k+1} = \text{prox}_{hf}(x^k)$

Es decir, si discretizamos la ec. dif. con Forward-Euler obtenemos el método de descenso por gradiente, y si discretizamos con Backward-Euler obtenemos el método proximal.

## Proximal Gradient Method

Consideremos ahora el problema  $\min_x f(x) + g(x)$   
con  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  y  $g: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  propias y convexas, y  $f$  diferenciable.

Usualmente se usa para partir el problema en una parte diferenciable, y otra que no necesariamente lo es.

El Proximal Gradient Method es:  $x^{k+1} = \text{prox}_{\lambda g}(x^k - \lambda \nabla f(x^k))$

Es decir, es como hacer un descenso por gradiente de la parte diferenciable, y luego aplicar el operador proximal de  $g$ .

Obs:

- Cuando  $g \equiv 0$ , obtenemos el clásico descenso por gradiente.
- Cuando  $g$  es la indicatriz de un conjunto convexo, obtenemos el Projected Gradient Descent.

### Interpretación

Como antes, estamos buscando puntos fijos. En este caso

$x^* = \text{prox}_{\lambda g}(x^* - \lambda \nabla f(x^*))$ . Veamos que este punto fijo es

la solución de  $\min_x f(x) + g(x)$ :

$$x^* = \text{prox}_{\lambda g}(x^* - \lambda \nabla f(x^*)) \Leftrightarrow x^* \text{ minimiza } g(x) + \frac{1}{2\lambda} \|x - x^* + \lambda \nabla f(x^*)\|_2^2$$

$$\Leftrightarrow 0 \in \partial g(x^*) + \frac{1}{\lambda} (x^* - x^* + \lambda \nabla f(x^*)) \Leftrightarrow 0 \in \partial g(x^*) + \nabla f(x^*)$$

$$\Leftrightarrow x^* \text{ minimiza } f(x) + g(x)$$

Obs En este caso, como hay un peso de gradiente,  $\lambda$  no puede ser arbitrario.

Por ejemplo si  $\nabla f$  es Lipschitz de constante  $L$ , entonces  $\lambda$  tiene

que estar en  $[0, \frac{2}{L})$ . Se puede hacer line-search también, o Armijo.