

Algoritmos de Optimización

No siempre se puede encontrar soluciones cerradas (casi nunca).

Se consideran entonces distintos algoritmos. Una clase grande corresponde a métodos iterativos.

- Ejemplos:
- "Bisección" (Golden section search)
 - Descenso (gradiente, Newton)
 - Proximal / Punto fijo
 - Punto interior
 - Stochastic GD

Métodos de descenso / de gradiente

Consideremos un problema de optimización sin restricciones, $\min_{x \in \mathbb{R}^n} f(x)$

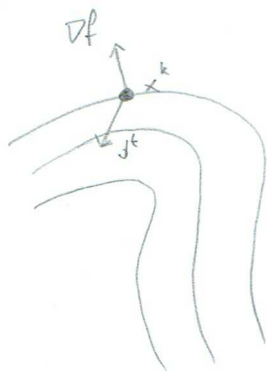
Empezamos de un punto x^0 , y generamos sucesivamente x^1, x^2, x^3, \dots

Si $f(x^{k+1}) < f(x^k) \quad \forall k$ decimos que es un método de descenso.

Tomemos $x^{k+1} = x^k + \alpha^k d^k$, donde $d^k \in \mathbb{R}^n$ es una dirección, α^k el tamaño del paso, positivo.

Un ejemplo puede ser $d^k = -\nabla f(x^k)$, o más en general,

$d^k / \langle \nabla f(x^k), d^k \rangle < 0$. A estos se los llama métodos de gradiente.



Esta condición garantiza que f decrece localmente en esa dirección. Es decir, $\exists \delta > 0 /$

$$f(x^k + \alpha^k d^k) < f(x^k) \quad \text{si } \alpha^k \in (0, \delta)$$

El método es entonces:

Dado un punto inicial x^0

Repetir

- || Elegir dirección d^k
- || Elegir peso α^k
- || Hacer $x^{k+1} = x^k + \alpha^k d^k$

Hasta condición de parada

¿Cómo elegir dirección d^k ?

¿Cómo elegir el peso α^k ?

¿Qué condición de parada usar?

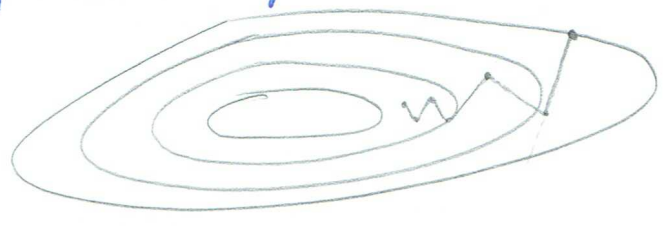
Elección de dirección

Escribamos penitenciate $x^{k+1} = x^k - \alpha^k D^k \nabla f(x^k)$ con $D^k \in S_{++}^n$

La condición $\nabla f(x^k)^T d^k < 0$ se traduce en $\nabla f(x^k)^T D^k \nabla f(x^k) > 0$ que se cumple pues D^k es definida positiva.

$D^k = I$

Corresponde a la dirección de máximo descenso, pero puede tener convergencia lenta, dependiendo de la geometría de las curvas de nivel.



Newton

$$D^k = (\nabla^2 f(x^k))^{-1}$$

Se le de considerar la aproximación de 2º orden: $f(x) \approx f(x^k) + \nabla f(x^k)(x-x^k) + \frac{1}{2}(x-x^k)^T \nabla^2 f(x^k)(x-x^k)$

$$\hookrightarrow x^{k+1} = x^k - (\nabla^2 f(x^k))^{-1} \nabla f(x^k) \quad (\text{corresponde a } \alpha^k = 1)$$

$$x^{k+1} = x^k - \alpha^k (\nabla^2 f(x^k))^{-1} \nabla f(x^k)$$

Diagonal Scaling

$$D^k = \begin{pmatrix} d_1^k & & 0 \\ & d_2^k & \\ 0 & & d_n^k \end{pmatrix}$$

$$\text{Cuando } d_i^k \equiv \left(\frac{\partial^2 f(x^k)}{\partial x_i^2} \right)^{-1}$$

es una aproximación a Newton.

Elección del paso α^k

Line search

Encontrar el α^k que minimiza f :

$$\alpha^k / f(x^k + \alpha^k d^k) = \min_{\alpha \geq 0} f(x^k + \alpha d^k)$$

Limited Line search

Se fija un $s > 0$ y se busca el mejor $\alpha^k \in [0, s]$

$$\alpha^k / f(x^k + \alpha^k d^k) = \min_{\alpha \in [0, s]} f(x^k + \alpha d^k)$$

Regla de Armijo

Encontrar el mejor α^k con line search suele ser costoso.

Una alternativa es probar con un paso inicial s , e ir reduciéndolo por un factor hasta encontrar un punto

donde f disminuye. Esto puede tener problemas de convergencia, pero con una modificación funciona bien.

Fijamos dos parámetros, $0 < \epsilon < 1$ y $0 < \beta < 1$.

Empezamos con $s (= 1)$ y vamos reduciendo por un factor β ($\beta^m s$) hasta que

$$f(x^k) - f(x^k + \beta^m s d^k) \geq -\epsilon \beta^m s \nabla f(x^k)^T d^k$$

Otra forma de verlo:

$$t = 1$$

Mientras

$$f(x^k + t d^k) > f(x^k) + \epsilon t \nabla f(x^k)^T d^k$$

$$t = \beta t$$

Interpretación (Armijo)

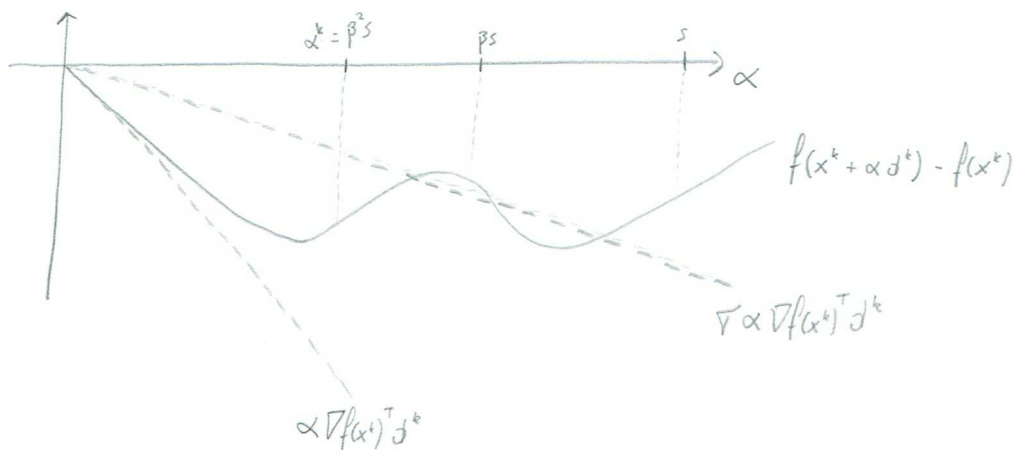
(10)

$$\underbrace{f(x^k) - f(x^k + \beta^m s d^k)}_{\text{ganancia al moverse a } x^k + \beta^m s d^k} \geq - \underbrace{\tau \beta^m s \nabla f(x^k)^T d^k}_{\text{ganancia si fuera lineal}}$$

nos conformamos con una fracción de la ganancia si fuera lineal

Usualmente: $\tau = 0.1$ ($\tau \in [10^{-5}, 10^{-1}]$)

$$\beta \in \left[\frac{1}{10}, \frac{1}{2} \right]$$



Peso constante

Se selecciona un peso $\alpha^k = s$ constante $\forall k$.

Si es muy grande, no se puede asegurar convergencia.

Si es muy chico, puede llevar a convergencia muy lenta.

(No se puede asegurar que sea de descenso en todo momento)

Peso decreciente

Hacemos $\alpha^k \rightarrow 0$, pero que $\sum \alpha^k = \infty$
(temporalmente se puede garantizar descenso en cada iteración)

Si $x^k \rightarrow \bar{x}$, por n, m grandes

$$x^n \approx x^m \approx \bar{x}$$

$$x^m \approx x^n - \left(\sum_{k=n}^{m-1} \alpha^k \right) \nabla f(\bar{x})$$

$$\Rightarrow \bar{x} \text{ estacionario } (\nabla f(\bar{x}) = 0)$$

Condición de parada

Como buscamos puntos críticos, una condición típica es

$$\|\nabla f(x^k)\| \leq \varepsilon$$

O similar. Por ejemplo $\frac{\|\nabla f(x^k)\|}{\|\nabla f(x^0)\|} \leq \varepsilon$ o $\|d^k\| \leq \varepsilon$

Análisis de convergencia

Decimos que las direcciones d^k son gradient related si para toda subsecuencia $\{x^{k_k}\}_{k \in K}$ que converge a un punto no estacionario, la subsecuencia correspondiente $\{d^{k_k}\}_{k \in K}$ es acotada y $\limsup_{k \rightarrow \infty, k \in K} \nabla f(x^{k_k})^T d^{k_k} < 0$

En particular, si $d^k = -D^k \nabla f(x^k)$, entonces basta ver que los valores propios de D^k estén alejados de 0: $c_1 \|z\|^2 \leq z^T D^k z \leq c_2 \|z\|^2 \quad \forall z \in \mathbb{R}^n$

Prop

Sea x^k generada por un método de gradiente $x^{k+1} = x^k + \alpha^k d^k$, con d^k gradient related y α^k según la regla de Armijo.

Entonces todo punto límite de x^k es un punto estacionario

\exists subsecuencia de x^k que converge a él

Dem

Sea $\{x^{k_k}\}_{k \in K}$ subsecuencia que converge a \bar{x} , y suponemos que $\nabla f(\bar{x}) \neq 0$.

$\{f(x^{k_k})\}$ es monótona decreciente, y $\{f(x^{k_k})\}_{k \in K} \rightarrow f(\bar{x})$ porque f es continua, por lo tanto $\{f(x^k)\} \rightarrow f(\bar{x})$ (toda la sucesión) $\Rightarrow f(x^k) - f(x^{k+1}) \rightarrow 0$

De la regla de Armijo: $f(x^k) - f(x^{k+1}) \geq -\sigma \alpha^k \nabla f(x^k)^T d^k \Rightarrow \alpha^k \nabla f(x^k)^T d^k \rightarrow 0$

Como d^k es gradient related, y suponemos \bar{x} no estacionario:

$$\limsup_{\substack{k \rightarrow \infty \\ k \in K}} \nabla f(x^{k_k})^T d^{k_k} < 0 \quad \Rightarrow \quad \alpha^k \rightarrow 0$$

Como en Armijo a cada iteración se empieza con paso s , si $\alpha^k \rightarrow 0$ entonces el paso inicial es reducido al menos una vez a partir de un \bar{k} ;

$$f(x^k) - f(x^k + \frac{\alpha^k}{\beta} d^k) < -\tau \frac{\alpha^k}{\beta} \nabla f(x^k)^T d^k \quad \forall k \geq \bar{k}$$

Le llamamos $\bar{\alpha}^k = \frac{\alpha^k}{\beta}$, que es el paso anterior al aceptado.

Reescribimos:

$$\frac{f(x^k) - f(x^k + \bar{\alpha}^k d^k)}{\bar{\alpha}^k} < -\tau \nabla f(x^k)^T d^k$$

Por Teorema de Valor Medio, $\exists \tilde{\alpha}^k \in [0, \bar{\alpha}^k]$ /

$$\frac{f(x^k) - f(x^k + \tilde{\alpha}^k d^k)}{\tilde{\alpha}^k} = -\nabla f(x^k + \tilde{\alpha}^k d^k)^T d^k$$

Entonces: $-\nabla f(x^k + \tilde{\alpha}^k d^k)^T d^k < -\tau \nabla f(x^k)^T d^k$

Como d^k es cotado (por ser gradient related) $\Rightarrow \exists d^{ke} \rightarrow \bar{d}$

Tomando limite según esta subsecuencia:

$$-\nabla f(\bar{x})^T \bar{d} < -\tau \nabla f(\bar{x})^T \bar{d}$$

$$\Rightarrow 0 \leq (1-\tau) \nabla f(\bar{x})^T \bar{d}$$

$$\text{y como } \tau < 1 \Rightarrow \nabla f(\bar{x})^T \bar{d} \geq 0$$

Lo cual es absurdo, pues contradice que d^k sea gradient related.