

IIMAS  
Universidad Nacional Autónoma de México

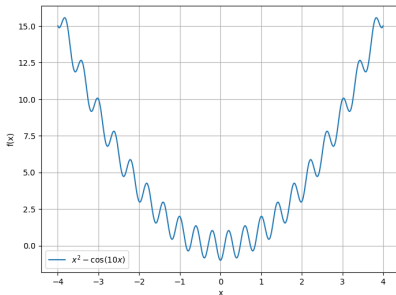
# **Agregando ruido a DG para llegar al mínimo global.**

Aprendizaje máquina teórico.

Alejandro Antonio Estrada Franco

# >Agregando ruido a DG para llegar al mínimo global

1



Para  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  el algoritmo que estudiaremos es el siguiente[1]:

**Algoritmo 1** : *Decenso de Gradiente Ruidoso*;  $x_1 \in \mathbb{R}^d$ ;  $\alpha > 0$ ;  $\sigma > 0$

- ▶ para  $t = 1, 2, \dots$
- ▶ Generar una realización de la v.a.  $U_t \sim \text{Unif}[B_{r=1}^2(0)]$
- ▶ Calcular  $\nabla f(x_t)$
- ▶  $x_{t+1} = x_t - \alpha \nabla f(x_t) + \sigma U_t$

# >Agregando ruido a DG para llegar al mínimo global

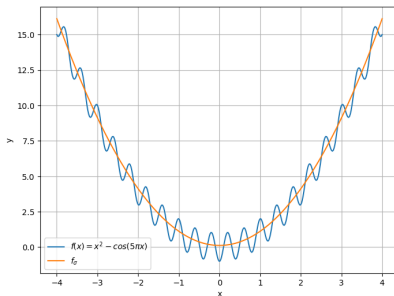
2

**Definición 1:** (componente  $\sigma$ -suave). Para todo  $\sigma > 0$  y una función  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , la componente  $\sigma$ -suave de  $f$  se define como:

$$f_\sigma(x) = \mathbb{E}_U[f(x + \sigma U)] = \int_{B_{r=1}^2(0)} f(x + \sigma y) p_U(y) dy$$

definimos también:

$$r_\sigma(x) = f(x) - f_\sigma(x)$$



# >Agregando ruido a DG para llegar al mínimo global

## Algunas observaciones y notación:

- Supondremos  $f$  continuamente diferenciable para obtener  $\nabla f_\sigma(x) = \nabla \mathbb{E}_U[f(x + \sigma U)] = \mathbb{E}_U[\nabla f(x + \sigma U)] = (\nabla f)_\sigma(x)$ ; es decir:

$$\nabla f_\sigma(x) = (\nabla f)_\sigma(x).$$

- También denotaremos  $\mathbb{E}_t[\cdot] = \mathbb{E}_{U_t}[\cdot]$ , así mismo denominaremos  $e_1 = x_1$ , y  $e_{t+1} = \mathbb{E}_t[x_{t+1}]$ , además  $\mathbb{E}[\cdot]$  denotará el valor esperado con respecto a la distribución conjunta de las uniformes (se toman independientes). Por ejemplo si  $x_t$  es el último punto del algoritmo, tenemos  $\mathbb{E}[x_t] = \mathbb{E}_1 \dots \mathbb{E}_{t-1}[x_t]$ .

# >Agregando ruido a DG para llegar al mínimo global

4

- ▶ Asumiremos también que la componente  $\sigma$ -suave  $f_\sigma$  es  $l$ -fuertemente convexa, y  $L$ -suave con  $0 < l \leq L < \infty$ . En consecuencia  $f_\sigma$  tiene un único minimizador  $x_\sigma^*$ , con  $f_\sigma^* = f_\sigma(x_\sigma^*)$ .
- ▶ Otro supuesto que se considera es que existen escalares  $M$  y  $\mu$  tales que los gradientes  $\nabla f$  y  $\nabla f_\sigma$  satisfacen:  
$$\mathbb{V}_U[\nabla f(e + \sigma U)] \leq M_\sigma + \mu \|\nabla f_\sigma(e)\|^2$$
- ▶ Finalmente consideraremos que  $f$  es  $L_0$ -Lipschitz con respecto a la norma euclidiana en  $\mathbb{R}^d$

# >Agregando ruido a DG para llegar al mínimo global

5

**Lema 1:** Bajo las iteraciones de DGR se satisface para todo  $n \in \mathbb{N}$ :  
 $\mathbf{e}_{t+1} = \mathbf{e}_t - \alpha \nabla f(\mathbf{e}_t + \sigma \mathbf{U}_{t-1}) + \sigma \mathbf{U}_{t-1}$  mas aún se cumple también  
 $\mathbb{E}_{t-1}[\mathbf{e}_{t+1}] = \mathbf{e}_t - \alpha \nabla f_\sigma(\mathbf{e}_t)$ .

*Demostración:* Para cada  $t \in \mathbb{N}$  tenemos  $\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha \nabla f(\mathbf{x}_t) + \sigma \mathbf{U}_t$ ,  
tomando esperanzas con respecto a  $\mathbf{U}_t$ ;  $\mathbf{e}_{t+1} = \mathbf{x}_{t+1} - \alpha \nabla f(\mathbf{x}_t)$  de  
donde  $\mathbf{e}_{t+1} = \mathbf{x}_{t+1} - \sigma \mathbf{U}_t$  luego  $\mathbf{e}_{t+1} - \mathbf{e}_t = -\alpha \nabla f(\mathbf{e}_t + \sigma \mathbf{U}_{t-1}) + \sigma \mathbf{U}_{t-1}$   
concluimos  $\mathbb{E}_{t-1}[\mathbf{e}_{t+1} - \mathbf{e}_t] = -\alpha(\nabla f)_\sigma(\mathbf{e}_t)$   $\square$ .

**Lema 2:** Se cumple  $f_\sigma(\mathbb{E}[\mathbf{x}_{t+1}]) \leq \mathbb{E}[f_\sigma(\mathbf{e}_{t+1})]$

*Demostración:* De la convexidad de  $f_\sigma$  y la desigualdad de Jensen:

$$\begin{aligned} f_\sigma(\mathbb{E}[\mathbf{x}_{t+1}]) &= f_\sigma(\mathbb{E}_1 \dots \mathbb{E}_{t-1} \mathbb{E}_t[\mathbf{x}_{t+1}]) \\ &= f_\sigma(\mathbb{E}_1 \dots \mathbb{E}_{t-1}[\mathbf{e}_{t+1}]) \\ &= f_\sigma(\mathbb{E}[\mathbf{e}_{t+1}]) \leq \mathbb{E}[f_\sigma(\mathbf{e}_{t+1})] \end{aligned}$$

$\square$ .

# >Agregando ruido a DG para llegar al mínimo global

**Lema 3:** Para todo  $t \in \mathbb{N}$  se verifica lo siguiente:

$$\mathbb{E}_{t-1} \|\mathbf{e}_{t-1} - \mathbf{e}_t\| \leq 2\alpha^2(\mu + 1) \|\nabla f_\sigma(\mathbf{e}_t)\|^2 + \alpha^2 M$$

donde:

$$M = 2M_\sigma + 2 \left( \frac{\sigma}{\alpha} \right)^2$$

**Teorema 1:** Para todo  $t \in \mathbb{N}$  se cumple:

$$f_\sigma(\mathbb{E}[x_{t+1}]) - f_\sigma^* - \frac{\alpha LM}{2l} \leq \rho^t \left( f_\sigma(x_1) - f_\sigma^* - \frac{\alpha LM}{2l} \right)$$

Si

$$0 < \alpha \leq \frac{1}{2L(\mu + 1)}$$

donde  $\rho = 1 - \alpha l \in (\frac{1}{2}, 1)$

# >Agregando ruido a DG para llegar al mínimo global

7

**Lema 4:** Para todo  $\sigma > 0$  se cumple  $|f_\sigma(x) - f(x)| \leq L_0\sigma$

*Demostración:*

$$\begin{aligned}|f_\sigma(x) - f(x)| &= |\mathbb{E}_U[f(x + \sigma U) - f(x)]| \\ &\leq \mathbb{E}_U|f(x + \sigma U) - f(x)| \\ &\leq L_0\mathbb{E}_U[|\sigma U|] = L_0\sigma\mathbb{E}_U[|U|] \leq L_0\sigma \quad \square.\end{aligned}$$

**Lema 5:**  $f_\sigma^* \leq f^* + L_0\sigma$

*Demostración:* Tenemos que  $f(x) = f_\sigma(x) + r_\sigma(x)$ , entonces:

$$f^* = \min_{x \in \mathbb{R}^d} f(x) \geq \min_{x \in \mathbb{R}^d} f_\sigma(x) + \min_{x \in \mathbb{R}^d} r_\sigma(x) = f_\sigma^* + \min_{x \in \mathbb{R}^d} r_\sigma(x)$$

Deacuerdo al Lema 4:

$$-L_0\sigma \leq \min_{x \in \mathbb{R}^d} r_\sigma \leq L_0\sigma$$

Entonces tenemos  $f^* \geq f_\sigma^* - L_0\sigma \quad \square.$



# >Agregando ruido a DG para llegar al mínimo global

**Teorema 2:** Para  $\alpha$  y  $\rho$  como en el Teorema 1 se cumple para todo  $t \in \mathbb{N}$ :

$$f(\mathbb{E}[x_{t+1}]) - f^* \leq \rho^t M_1 + M_2$$

Donde:

$$M_1 = f_\sigma(x_1) - f_\sigma^* - \frac{\alpha L (M_\sigma + (\frac{\sigma}{\alpha})^2)}{l},$$
$$M_2 = \alpha \left( \frac{LM_\sigma}{l} + \frac{L(\frac{\sigma}{\alpha})^2}{l} + L_0 \left( \frac{\sigma}{\alpha} \right) \right) + L_0 \sigma$$

*Demostración:* Por el Lema 4 tenemos para todo  $t \in \mathbb{N}$ ;  
 $f(\mathbb{E}[x_{t+1}]) \leq f_\sigma(\mathbb{E}[x_{t+1}]) + L_0 \sigma = f_\sigma(\mathbb{E}[x_{t+1}]) + \alpha L_0 \left( \frac{\sigma}{\alpha} \right)$ , luego del Lema 5  $-f^* \leq -f_\sigma^* + L_0 \sigma$  por lo que tenemos:

$$f(\mathbb{E}[x_{t+1}]) - f^* \leq f_\sigma(\mathbb{E}[x_{t+1}]) - f_\sigma^* + \alpha L_0 \left( \frac{\sigma}{\alpha} \right) + L_0 \sigma$$

# >Agregando ruido a DG para llegar al mínimo global

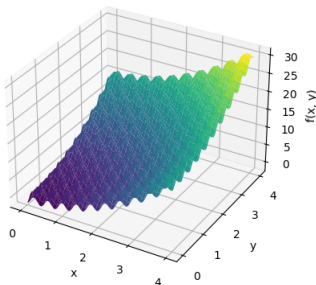
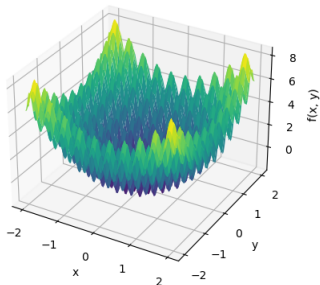
Del Teorema 1 tenemos:

$$\begin{aligned} f(\mathbb{E}[x_{t+1}]) - f^* &\leq f_\sigma(\mathbb{E}[x_{t+1}]) - f_\sigma^* + \alpha L_0 \left( \frac{\sigma}{\alpha} \right) + L_0 \sigma \\ &\leq \rho^t \left( f_\sigma(x_1) - f_\sigma^* - \frac{\alpha LM}{2l} \right) + \frac{\alpha LM}{2l} + \alpha L_0 \left( \frac{\sigma}{\alpha} \right) + L_0 \sigma \\ &= \rho^t M_1 + \alpha \left( \frac{2M_\sigma + 2 \left( \frac{\sigma}{\alpha} \right)^2}{2l} + L_0 \left( \frac{\sigma}{\alpha} \right) \right) + L_0 \sigma \\ &= \rho^t M_1 + M_2 \quad \square. \end{aligned}$$

# >Agregando ruido a DG para llegar al mínimo global

10

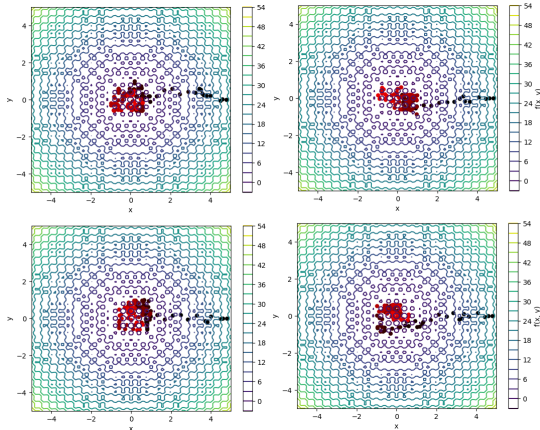
**Ejemplo:** Consideremos la función  
 $f(x, y) = x^2 + y^2 - \cos(5\pi x) - \cos(5\pi y)$



# >Agregando ruido a DG para llegar al mínimo global

11

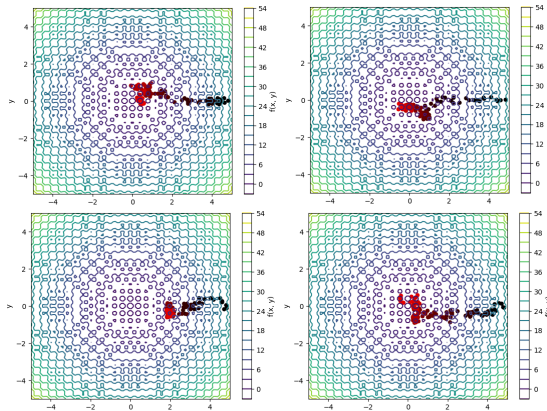
Se ilustran cuatro ejecuciones del algoritmo GDR con  $t = 100$ ,  $x_1 = (5, 0)$ ,  $\sigma = 0,05$ ,  $\alpha = 0,02$



# >Agregando ruido a DG para llegar al mínimo global

12

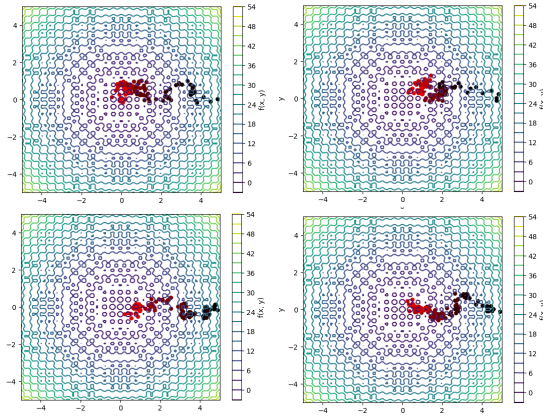
Se ilustran cuatro ejecuciones del algoritmo GDR con  $t = 100$ ,  $x_1 = (5, 0)$ ,  $\sigma = 0,01$ ,  $\alpha = 0,015$



# > Agregando ruido a DG para llegar al mínimo global

13

Se ilustran cuatro ejecuciones del algoritmo GDR con  $t = 100$ ,  $x_1 = (5, 0)$ ,  $\sigma = 0,015$ ,  $\alpha = 0,015$



# >Agregando ruido a DG para llegar al mínimo global

- [1] Xuliang Quin, Xin Xiu y Xiaopeng Luo. «Global Convergence of Noisy Gradient Descent». En: *IEEE International Conference on Systems, Man, and Cybernetics (SMC)* (2022).