

Non linear optimizers

Combinatorial Decision Making and Optimization Project

Marini Luca

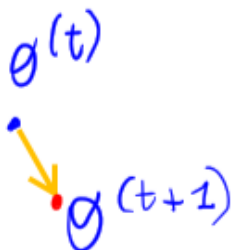
1. Optimizers

2. Tested functions

Optimizers

$$\theta^{(t+1)} = \theta^t - lr \nabla f(\theta^{(t)})$$

Gradient descent



\times

\times = global optimum

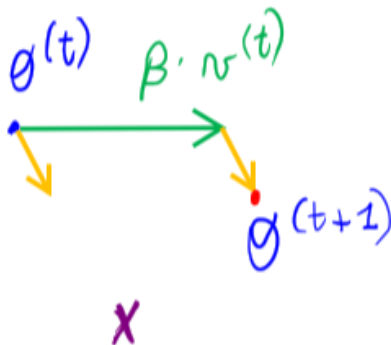
$$\rightarrow = -\eta \cdot \nabla f(\theta^{(t)})$$

$$v^{(t+1)} = \beta v^{(t)} - lr \nabla f(\theta^{(t)})$$

$$\theta^{(t+1)} = \theta^{(t)} + v^{(t+1)}, \text{ with } \beta \in [0, 1)$$

if $\beta = 0$, then the update step is the same as the one of gradient descent.

Momentum

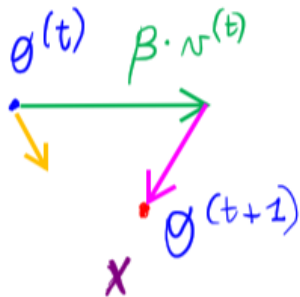


X = global optimum

\rightarrow = $-\eta_r \cdot \nabla f(\theta^{(t)})$

$$\begin{aligned}v^{(t+1)} &= \beta v^{(t)} - lr \nabla f(\theta^{(t)} + \beta v^{(t)}) \\ \theta^{(t+1)} &= \theta^{(t)} + v^{(t+1)}, \text{ with } \beta \in [0, 1)\end{aligned}$$

Nesterov momentum



\times = global optimum

$$\text{magenta arrow} = -\eta_r \cdot \nabla f(\theta^{(t)} + \beta \cdot v^{(t)})$$

$$\text{yellow arrow} = -\eta_r \cdot \nabla f(\theta^{(t)})$$

Adaptive learning rates method

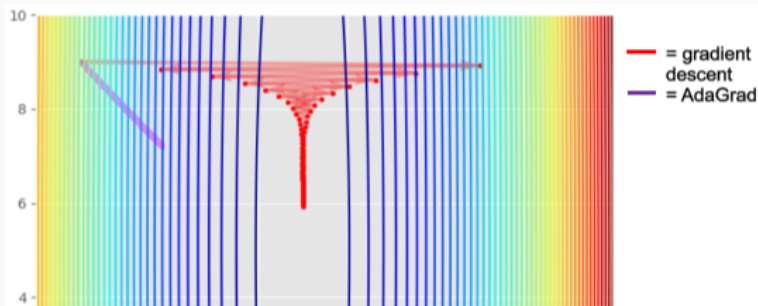
- Methods that adapt the learning rate to the parameters
- they reduce the step of updates for dimensions whose gradient direction is not consistent across iterations
- they increase the step of updates for dimensions whose gradient direction is consistent across iterations

AdaGrad (Adaptive Gradient)

$$\begin{aligned}s^{(t+1)} &= s^{(t)} + \nabla f(\theta^{(t)}) \nabla f(\theta^{(t)}) \\ \theta^{(t+1)} &= \theta^{(t)} - \frac{lr}{\sqrt{s^{(t+1)}} + \epsilon} \nabla f(\theta^{(t)})\end{aligned}$$

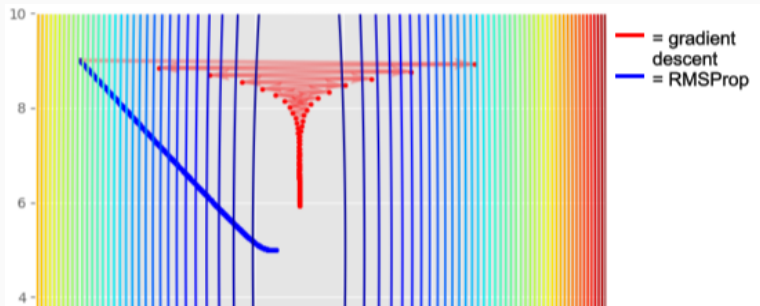
AdaGrad (Adaptive Gradient)

$$f(x_1, x_2) = 90(x_1 - 3)^2 + (x_2 - 5)^2$$



$$s^{(t+1)} = \beta s^{(t)} + (1 - \beta) \nabla f(\theta^{(t)}) \nabla f(\theta^{(t)})$$
$$\theta^{(t+1)} = \theta^{(t)} - \frac{lr}{\sqrt{s^{(t+1)} + \epsilon}} \nabla f(\theta^{(t)}) \text{ , with } \beta \geq 0.9$$

$$f(x_1, x_2) = 90(x_1 - 3)^2 + (x_2 - 5)^2$$



$$g^{(t+1)} = \beta_1 g^{(t)} + (1 - \beta_1) \nabla f(\theta^{(t)})$$

$$s^{(t+1)} = \beta_2 s^{(t)} + (1 - \beta_2) \nabla f(\theta^{(t)}) \nabla f(\theta^{(t)})$$

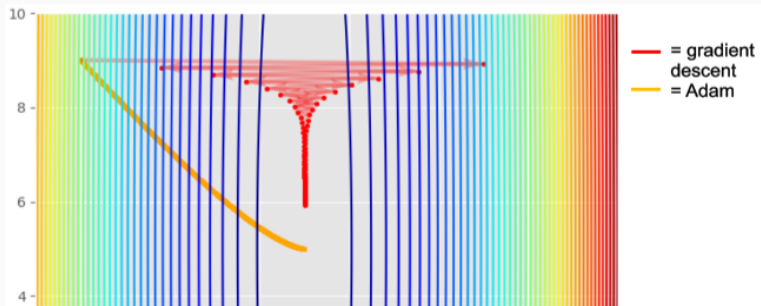
$$g^{debiased} = \frac{g^{(t+1)}}{1 - \beta_1^{t+1}}$$

$$s^{debiased} = \frac{s^{(t+1)}}{1 - \beta_2^{t+1}}$$

$$\theta^{(t+1)} = \theta^{(t)} - \frac{lr}{\sqrt{s^{debiased} + \epsilon}} g^{debiased}$$

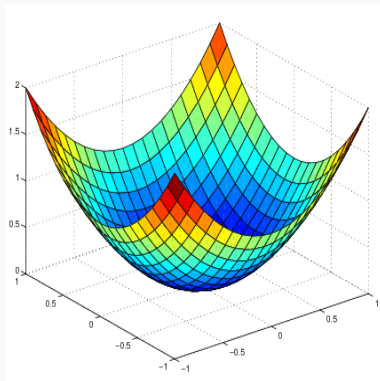
$$\beta_1 = 0.9, \beta_2 = 0.999.$$

$$f(x_1, x_2) = 90(x_1 - 3)^2 + (x_2 - 5)^2$$



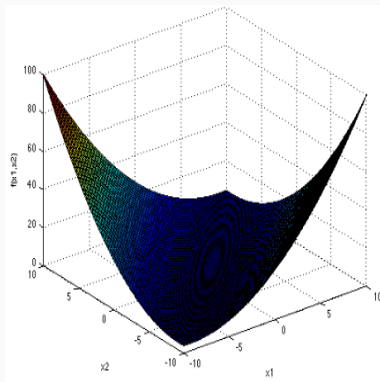
Tested functions

Paraboloid



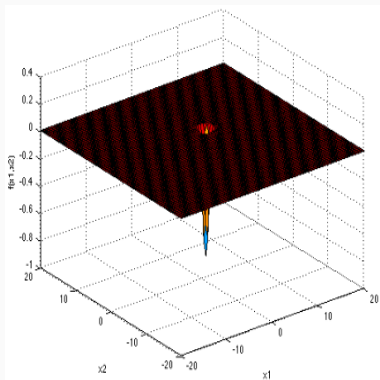
- Function: $x_1^2 + x_2^2$
- Global Minimum: $x = (0, 0)$

Matyas function



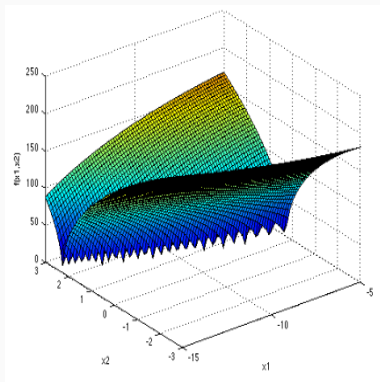
- Function: $0.26(x_1^2 + x_2^2) - 0.48x_1x_2$
- Global Minimum: $x = (0, 0)$

Easom function



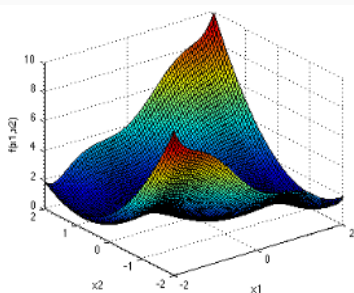
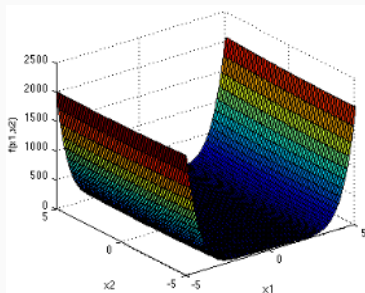
- Function: $-\cos(x_1)\cos(x_2)\exp(-(x_1 - \pi)^2 - (x_2 - \pi)^2)$
- Global Minimum: $x = (\pi, \pi)$

Bukin function



- Function: $100\sqrt{|x_2 - 0.01x_1^2|} + 0.01|x_1 + 10|$
- Global Minimum: $x = (-10, 1)$

Three hump camel function



- Function: $2x_1^2 - 1.05x_1^4 + \frac{x_1^6}{6} + x_1x_2 + x_2^2$
- Global Minimum: $x = (0, 0)$