

FGV EMAp  
João Pedro Jerônimo

Otimização para Ciência de Dados  
Revisão para A2

Rio de Janeiro  
2025

# Conteúdo

- 1 Introdução ..... 3
- 2 Método do Gradiente ..... 5
  - 2.1 Caso Global ..... 6
  - 2.2 Caso Convexo ..... 8
  - 2.3 Interpretação via regularização ..... 9
- 3 Método do Subgradiente ..... 10

# Introdução

Antes de iniciarmos com o conteúdo de verdade, vou relembrar alguns conceitos importantes da A1 que vão ajudar a entender os métodos presentes nesse PDF.

**Teorema 1.1** (Aproximação de Primeira Ordem): Quando  $f$  é continuamente diferenciável, em uma vizinhança de um ponto  $x$  podemos mostrar que:

$$\forall d \in \mathbb{R}^n \text{ com } \|d\| = 1 \quad \frac{df}{dd} = \nabla f(x)^T d \quad (1)$$

e, além disso, temos:

$$\forall y \in \mathbb{R}^n \text{ na vizinhança} \quad f(y) = f(x) + \nabla f(x)^T (y - x) + o(\|y - x\|) \quad (2)$$

Onde  $o : \mathbb{R}_+ \rightarrow \mathbb{R}$  satisfaz  $\lim_{t \rightarrow 0^+} \frac{o(t)}{t} = 0$

**Teorema 1.2** (Aproximação Linear): Seja  $f : U \rightarrow \mathbb{R}$  uma função duas vezes continuamente diferenciável e  $U \subseteq \mathbb{R}^n$ , e seja  $x \in U$  e  $r > 0$  tais que  $B(x, r) \subset U$  então:

$$\begin{aligned} \forall y \in B(x, r) \quad \exists \xi \in [x, y] \text{ tal que} \\ f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(\xi) (y - x) \end{aligned} \quad (3)$$

**Teorema 1.3** (Aproximação de Segunda Ordem): Seja  $f : U \rightarrow \mathbb{R}$  uma função duas vezes continuamente diferenciável e  $U \subseteq \mathbb{R}^n$ , e seja  $x \in U$  e  $r > 0$  tais que  $B(x, r) \subset U$  então:

$$\begin{aligned} \forall y \in B(x, r) \text{ vale} \\ f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(x) (y - x) + o(\|y - x\|^2) \end{aligned} \quad (4)$$

# **Método do Gradiente**

## 2.1 Caso Global

É o método de otimização mais clássico que existe! Vamos supor que queremos resolver o problema:

$$\min_{x \in \mathbb{R}^n} f(x) \quad (5)$$

Lembra do que vimos em cálculo? Que  $\nabla f(x)$  é o vetor que aponta pra direção em que  $f(x)$  aumenta? Então que tal a gente seguir na direção contrária a  $f(x)$ ? Isso faz bastante sentido, e funciona! Mas deve ter um motivo mais matemático por trás, não é? Vamos primeiro mostrar o algoritmo:

---

```
1 func GradientDescent( $f$ ) {  
2    $x^{(0)} \in \mathbb{R}^n$   
3    $\alpha > 0$   
4   for  $t \in [T]$  do {  
5      $x^{(t+1)} = x^{(t)} - \alpha \nabla f(x^{(t)})$   
6   }  
7   return  $x^{(T)}$   
8 }
```

---

Algoritmo 1: Gradient Descent

Antes de entender um motivo mais matemático por trás do algoritmo, vamos ver algumas definições

**Definição 2.1.1** (Funções  $M$ -Lipschitz): Dizemos que uma função  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  é  $M$ -Lipschitz quando:

$$\|f(x) - f(y)\| \leq M \|x - y\| \quad \forall x, y \in \mathbb{R}^n \quad (6)$$

**Definição 2.1.2** (Função  $L$ -Suave): Uma função  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  é  $L$ -suave quando seu gradiente é  $L$ -Lipschitz:

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad \forall x, y \in \mathbb{R}^n \quad (7)$$

Essa definição de suavidade tem uma interpretação, imagine que, se eu estou na posição  $x$  e vou pra posição  $y$ , a variação que eu vou ter na função, dentro dessa passada, não ultrapassa o quanto eu andei vezes uma constante  $L$ . Então funções muito onduladas, e com ondulações

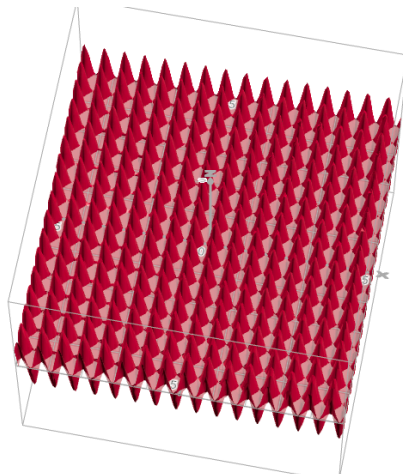


Figura 2: Função bem desregular, mas suave,  $f(x) = \sin(10x) + \cos(10y)$

**Teorema 2.1.1** (Aproximação linear de funções suaves): Seja  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  uma função diferenciável. Então  $f$  é  $L$ -suave se, e somente se,  $\forall x, y \in \mathbb{R}^n$ :

$$|f(y) - f(x) + \nabla f(x)^T(y - x)| \leq \frac{L}{2} \|y - x\|_2^2 \quad (8)$$

Usando esse teorema, a gente pode escrever isso:

$$\begin{aligned} f(x^{(t+1)}) &\leq f(x^{(t)}) + \nabla f(x^{(t)})^T(x^{(t+1)} - x^{(t)}) + \frac{L}{2} \|x^{(t+1)} - x^{(t)}\|^2 \\ &\leq f(x^{(t)}) - \alpha \|\nabla f(x^{(t)})\|^2 + \frac{\alpha^2 L}{2} \|\nabla f(x^{(t)})\|^2 \\ &\leq f(x^{(t)}) - \alpha \left(1 - \frac{\alpha L}{2}\right) \|\nabla f(x^{(t)})\|^2 \end{aligned} \quad (9)$$

E isso vale quando  $\alpha \in (0, \frac{2}{L})$ , ou seja, se o ponto atual  $x^{(t)}$  **NÃO É ESTACIONÁRIO**, o valor da função no próximo ponto será **menor** que o valor do ponto atual menos o tamanho do gradiente ao quadrado vezes um termo de regulação. Parece complicado, mas o que isso quer dizer? Eu vou usar esse fato para mostrar que, independente do ponto que eu iniciar o método do gradiente, eu **sempre vou encontrar um ponto mínimo local utilizando o método do gradiente**

**Teorema 2.1.2** (Convergência do Gradient Descent): Suponha que  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  é  $L$ -suave. Tome qualquer passo:

$$\alpha = \frac{\beta}{L} \quad (10)$$

para algum  $\beta \in (0, 2)$ . Então:

$$\min_{t \in [T]} \|\nabla f(x^{(t)})\|_2^2 \leq \frac{1}{T} \sum_{t=1}^T \|\nabla f(x^{(t)})\|_2^2 \leq \left(\frac{2/\beta}{2-\beta}\right) \frac{L(f(x^{(1)}) - f^*)}{T} \quad (11)$$

*Demonstração:* Sabemos que, dado  $n$  pontos  $x_i$ , a média  $\frac{1}{n} \sum_{i=1}^n x_i \in [\min(x_i), \max(x_i)]$ , então a desigualdade inicial já está provada. Vamos provar a segunda. Usando a equação (9), temos que:

$$\alpha \left(1 - \frac{\alpha L}{2}\right) \|\nabla f(x^{(t)})\|^2 \leq f(x^{(t)}) - f(x^{(t+1)}) \quad (12)$$

isso para todo  $t \in [T]$ , então vamos somar todos os termos para obter:

$$\begin{aligned} \alpha \left(1 - \frac{\alpha L}{2}\right) \sum_{t=1}^T \|\nabla f(x^{(t)})\|_2^2 &\leq \sum_{t=1}^T (f(x^{(t)}) - f(x^{(t+1)})) \\ &\leq f(x^{(1)}) - f(x^{(T+1)}) \\ &= f(x^{(1)}) - f^* + f^* + f(x^{(T+1)}) \\ &\leq f(x^{(1)}) - f^* \end{aligned} \quad (13)$$

A primeira desigualdade eu fiz uma soma telescópica, depois eu somei 0 ( $f^* - f^*$ ) e, como  $f^*$  é o valor mínimo da função, com certeza subtrair a parte que eu somei  $f^*$  vai dar um valor maior, então eu obtenho o resultado do enunciado do teorema dividindo tudo por  $\alpha(1 - \frac{\alpha L}{2})T$   $\square$

Por que esse teorema mostra que, independentemente do lugar, o algoritmo converge para um ponto estacionário? Ele ta me dizendo isso daqui:

$$\min_{t \in [T]} \|\nabla f(x^{(t)})\|_2^2 = O\left(\frac{1}{T}\right) \quad (14)$$

Ou seja, o mínimo **converge para 0** conforme  $T \rightarrow \infty$  **independentemente do ponto inicial**  $x^{(0)}$

Só que se pararmos para pensar, se  $\nabla f(x)$  é uma direção de subida, então  $-\nabla f(x)$  é de descida, mas será que é a melhor? Será que **precisa** ser a melhor para o método funcionar? Se eu escolher **uma direção de descida arbitrária** ele funciona? Mas antes disso, vamos entender o que é uma direção de descida

**Definição 2.1.3** (Direção de Descida): Dizemos que  $d \in \mathbb{R}^n$  é uma direção de descida para a função  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  no ponto  $x \in \mathbb{R}^n$  quando:

$$d^T \nabla f(x) < 0 \quad (15)$$

então podemos considerar um novo algoritmo

---

```

1 func GradientDescentWithArbitraryDescentDirection( $f$ ) {
2    $x^{(0)} \in \mathbb{R}^n$ 
3    $\alpha > 0$ 
4   for  $t \in [T]$  do {
5     Encontrar uma direção de descida  $d^{(t)}$ 
6      $x^{(t+1)} = x^{(t)} + \alpha d^{(t)}$ 
7   }
8   return  $x^{(T)}$ 
9 }
```

---

Algoritmo 3: Gradient Descent com Direção de Descida Arbitrária

Podemos provar, analogamente ao Teorema 2.1.2, que o algoritmo converge para um ponto estacionário

## 2.2 Caso Convexo

Antes, não assumimos nada além da suavidade da função, agora vamos mostrar que, assumindo que  $f$  é convexa, o algoritmo converge para uma solução global. Primeiro, nós sabemos que:

$$\begin{aligned} x^{(t+1)} &= x^{(t)} - \alpha \nabla f(x^{(t)}) \\ \Leftrightarrow x^{(t+1)} - x^* &= x^{(t)} - x^* - \alpha \nabla f(x^{(t)}) \\ \Leftrightarrow \|x^{(t+1)} - x^*\|_2^2 &= \|x^{(t)} - x^* - \alpha \nabla f(x^{(t)})\|_2^2 \end{aligned} \quad (16)$$

Pelas propriedades da convexidade:

$$(x^* - x^{(t)})^T \nabla f(x^{(t)}) \leq f(x^*) - f(x^{(t)}) \quad (17)$$

e pelo que vimos na equação (9), se  $\alpha \in (0, \frac{2}{L})$ , podemos chegar que:



$$\|x^{(t+1)} - x^*\|_2^2 \leq \|x^{(t)} - x^*\|_2^2 - \alpha \left( 2 - \frac{1}{1 - \frac{\alpha L}{2}} \right) (f(x^{(t)}) - f^*) \quad (18)$$

ou seja, a distância do próximo iterado pro ponto ótimo é menor a distância atual, menos um termo proporcional a distância dos resultados de  $x^{(t)}$  e do ponto ótimo. Vamos usar isso para provar a convergência global do resultado

**Teorema 2.2.1** (Convergência Convexa do Método do Gradiente): Suponha que  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  é  $L$ -suave e convexa e tome qualquer passo

$$\alpha = \frac{\beta}{L} \quad (19)$$

para algum  $\beta \in (0, 1)$ , então:

$$f(x^{(T)}) - f^* \leq \frac{1}{T} \sum_{t=1}^T (f(x^{(t)}) - f^*) \leq \frac{\beta^{-1} - \frac{1}{2}}{1 - \beta} \frac{L \|x^{(1)} - x^*\|_2^2}{T} \quad (20)$$

*Demonstração:* Somando-se em  $\text{tin}[T]$  a recorrência:

$$\alpha \left( 2 - \frac{1}{1 - \alpha \frac{L}{2}} \right) (f(x^{(t)}) - f^*) \leq \|x^{(t)} - x^*\|_2^2 - \|x^{(t+1)} - x^*\|_2^2 \quad (21)$$

obtemos novamente uma soma telescópica:

$$\begin{aligned} \alpha \left( \frac{1 - \alpha L}{1 - \alpha \frac{L}{2}} \right) \sum_{t=1}^T (f(x^{(t)}) - f^*) &\leq \sum_{t=1}^T (\|x^{(t)} - x^*\|_2^2 - \|x^{(t+1)} - x^*\|_2^2) \\ &\leq \|x^{(1)} - x^*\|_2^2 - \|x^{(T+1)} - x^*\|_2^2 \\ &\leq \|x^1 - x^*\|_2^2. \end{aligned} \quad (22)$$

Dividindo-se por  $\alpha \left( \frac{1 - \alpha L}{1 - \alpha \frac{L}{2}} \right) T$ , e lembrando que, pelo , a sequência  $\{f(x^{(t)})\}$  é decrescente:

$$f(x^{(T)}) - f^* \leq \frac{1}{T} \sum_{\{t=1\}}^T (f(x^{(t)}) - f^*) \leq \frac{\alpha^{-1} (1 - \alpha \frac{L}{2})}{T} \|x^1 - x^*\|_2^2. \quad (23)$$

□

## 2.3 Interpretação via regularização

Outra formas que podemos ver e interpretar o algoritmo do gradiente é resolver a seguinte fórmula:

$$x^{(t+1)} = \underset{x \in \mathbb{R}^n}{\operatorname{argmin}} \left( \underbrace{f(x^{(t)}) + (x - x^{(t)})^T \nabla f(x^{(t)})}_{\text{Aproximação Linear}} + \underbrace{\frac{1}{2\alpha^{(t)}} \|x - x^{(t)}\|_2^2}_{\text{Regularização Proximal}} \right) \quad (24)$$

Ou seja, eu vou pegar qual que é o valor que minimiza a aproximação linear regularizada por um termo quadrático

## **Método do Subgradiente**

Até agora vimos um método que usa o Gradiente da função, logo, para que ele funcione, estamos assumindo que a função é diferenciável em todos os pontos, porém em aplicações reais muitas funções não são diferenciáveis em todos os pontos. Então faz sentido utilizar esse método? Claro que não, porém, podemos utilizar uma versão muito parecida!

**Definição 3.1** (Subgradiente): Seja uma função  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  e dado  $x \in \mathbb{R}^n$ , um vetor  $g_x \in \mathbb{R}^n$  é chamado de **subgradiente** de  $f$  em  $x$  quando

$$f(y) \geq f(x) + (y - x)^T g_x \quad \forall y \in \mathbb{R}^n \quad (25)$$

**Definição 3.2** (Subdiferencial): O conjunto de todos os subgradientes de  $f$  em  $x$  é chamado de **subdiferencial** de  $f$  em  $x$  (Denotado por  $\partial f(x)$ )

$$\partial f(x) := \{g_x \in \mathbb{R}^n : f(y) \geq f(x) + (y - x)^T g_x \quad \forall y \in \mathbb{R}^n\} \quad (26)$$

O que seria um subgradiente **intuitivamente** então? Note que, se eu definir  $y = x^*$  (Sendo o ponto mínimo), vamos obter o seguinte:

$$f(x^*) \geq f(x) + g^T(x^* - x) \quad (27)$$

sabendo que  $f(x^*) \leq f(x)$ , temos que:

$$f(x) \geq f(x) + g^T(x^* - x) \Leftrightarrow 0 \geq g^T(x^* - x) \Leftrightarrow 0 \leq g^T(x - x^*) \quad (28)$$

Ou seja, o subgradiente faz um ângulo **maior que 90°** com o vetor que aponta de  $x$  para  $x^*$ . O que isso quer dizer? Que os subgradientes são direções que, se seguirmos na **direção oposta**, nós **não** estamos indo em uma direção em que nós temos **certeza** que ela sobe

Nem sempre existirão subgradientes, porém, quando falamos das **funções convexas**, mesmo elas não sendo **diferenciáveis**, elas possuem um subgradiente. Logo, um subgradiente é uma direção que, se eu vou na direção contrária a ela, eu tenho **certeza** que minha função não está aumentando

**Teorema 3.1:** Seja  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , então vale que:

$$f \text{ é convexa} \Leftrightarrow \partial f(x) \neq \emptyset \quad \forall x \quad (29)$$

Vale ressaltar que funções subdiferenciáveis podem não ser suaves. E por que isso é importante? Acontece que antes, no método do gradiente e na direções de descida, utilizamos o fato das funções serem suaves para provar a convergência do método, porém, aqui estamos trabalhando com funções que não necessariamente são diferenciáveis, logo, intuitivamente, elas podem ter vários picos, ou paradas bruscas, etc.

---

```

1 func SubgradientMethod( $f$ ) {
2    $x^{(1)} \in \mathbb{R}^n$ 
3    $\{\alpha^{(t)}\} \subset (0, \infty)$ 
4   for  $t \in [T]$  do {
5     Compute um subgradiente  $g^{(t)}$  de  $f$  em  $x^{(t)}$ 
6      $x^{(t+1)} = x^{(t)} - \alpha^{(t)} g^{(t)}$ 
7   }
8   return  $x^{(T)} := \frac{1}{T} \sum_{t=1}^T \frac{\alpha^{(t)}}{\sum_{l=1}^T \alpha^{(l)}} x^{(t)}$ 
9 }

```

---

Algoritmo 4: Método do Subgradiente

Esse algoritmo parece até que “ingênuo”, tipo, nada garante que o subgradiente vai fazer com que a função desça né?? Vamos mostrar que na verdade esse método converge sim! Porém, vamos assumir algumas coisas também. Para esse caso, vamos assumir que a função é  $M$ -Lipschitz (Definição 2.1.1)

Usando o que foi mostrado na introdução (Teorema 1.2), podemos mostrar o seguinte teorema:

**Teorema 3.2** (Aproximação Linear de funções  $M$ -Lipschitz): Seja  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  uma função convexa, então  $f$  é  $M$ -Lipschitz contínua se, e somente se,  $\forall x, y \in \mathbb{R}^n$  e todo subgradiente  $g_x \in \mathbb{R}^n$  de  $f$  em  $x$ :

$$|f(y) - f(x) + g_x^T(y - x)| \leq M \|y - x\|_2^2 \quad (30)$$

E no que isso me é útil? Só parece um bando de complicação esquisita. Na verdade, o que será útil na demonstração é uma **consequência** desse teorema

**Corolário 3.2.1:** Seja  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  uma função convexa, então  $f$  é  $M$ -Lipschitz se, e somente se,  $\forall x \in \mathbb{R}^n$  e todo subgradiente  $g_x$  de  $f$  em  $x$ , vale que  $\|g_x\|_2 \leq M$

Vamos então mostrar a convergência do algoritmo. Um ponto interessante a se dizer é que, como você **talvez** possa ter pensado, nem sempre um subgradiente vai ser uma direção de **descida**. O que isso quer dizer? Isso quer dizer que, não necessariamente, a cada iteração,  $f(x^{(t+1)}) \leq f(x^{(t)})$ , porém, ele decresce o valor sobre a **média de todas iterações**.

**Teorema 3.3** (Convergência do Subgradiente): Suponha que  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  é  $M$ -Lipschitz contínua e convexa. Então:

$$f(\bar{x}^{(T)}) - f^* \leq \frac{1}{T} \sum_{t=1}^T (f(x^{(t)}) - f^*) \leq \frac{\|x^{(1)} - x^*\|_2^2 + M^2 \sum_{t=1}^T (\alpha^{(t)})^2}{\sum_{t=1}^T \alpha^{(t)}} \quad (31)$$

em particular, para qualquer  $\beta > 0$ , consideramos:

$$\alpha^{(t)} = \frac{\beta}{\sqrt{T}} \quad t \in [T] \quad (32)$$

*Demonstração:* Primeiramente, temos que:

$$\begin{aligned}
\|x^{(t+1)} - x^*\|_2^2 &= \|x^{(t)} - x^* - \alpha^{(t)} g^{(t)}\|_2^2 \\
&\leq \|x^{(t)} - x^*\|_2^2 + 2\alpha^{(t)}(x - x^*)^T g^{(t)} + (\alpha^{(t)})^2 \|g^{(t)}\|_2^2
\end{aligned} \tag{33}$$

e pela definição de subgradiente ( $f$ , por ser convexa, é garantida de ter subgradientes pelo Teorema 3.1), temos que:

$$(x^* - x^{(t)})^T g^{(t)} \leq f(x^*) - f(x^{(t)}) \tag{34}$$

A partir disso, também podemos escrever:

$$\|x^{(t+1)} - x^*\|_2^2 \leq \|x^{(t)} - x^*\|_2^2 - \alpha^{(t)}(f(x^{(t)}) - f^*) + (M\alpha^{(t)})^2 \tag{35}$$

Agora nós vamos fazer novamente a soma por recursão:

$$\begin{aligned}
\sum_{t=1}^T \alpha^{(t)}(f(x^{(t)}) - f^*) &\leq \sum_{t=1}^T (\|x^{(t)} - x^*\|_2^2 - \|x^{(t+1)} - x^*\|_2^2) + M^2 \sum_{t=1}^T (\alpha^{(t)})^2 \\
&\leq \|x^{(1)} - x^*\|_2^2 + M^2 \sum_{t=1}^T (\alpha^{(t)})^2
\end{aligned} \tag{36}$$

Dividindo por  $\sum_{t=1}^T \alpha^{(t)}$  e, usando a desigualdade de Jensen

$$f(\bar{x}^{(T)}) - f^* \leq \frac{1}{T} \sum_{t=1}^T (f(x^{(t)}) - f^*) \leq \frac{\|x^{(1)} - x^*\|_2^2 + M^2 \sum_{t=1}^T (\alpha^{(t)})^2}{\sum_{t=1}^T \alpha^{(t)}} \tag{37}$$

Ou seja, a média das iterações converge para próximo de  $x^*$  □

Nós assumimos que  $\alpha^{(t)}$  está em um conjunto de passos, e não que é um único passo, mas e se assumirmos que é um único, qual seria o melhor passo? Tomando  $\alpha^{(t)} = \alpha$ :

$$f(\bar{x}^{(T)}) - f^* \leq \frac{\|x^{(1)} - x^*\|_2^2}{T\alpha} + M^2\alpha \tag{38}$$

perceba também que:

$$\min_{\alpha>0} \left\{ \frac{\|x^{(1)} - x^*\|_2^2}{T\alpha} + M^2\alpha \right\} = \frac{M \|x^{(1)} - x^*\|_2}{\sqrt{T}} \tag{39}$$

logo, temos que:

$$\alpha^{(t)} = \frac{\|x^{(1)} - x^*\|_2}{M\sqrt{T}} \tag{40}$$

então teremos a taxa de convergência:

$$f(\bar{x}^{(T)}) - f^* \leq \frac{2M \|x^{(1)} - x^*\|_2}{\sqrt{T}} \tag{41}$$

Porém, na maioria esmagadora das vezes, não sabemos  $M$  e  $\|x^{(1)} - x^*\|$ , então utilizamos o passo sequencial já definido anteriormente