

FGV EMap
João Pedro Jerônimo

Otimização para Ciência de Dados

Revisão para A1

Rio de Janeiro
2025

Conteúdo

1	Otimização Irrestrita	3
1.1	Introdução	4
1.2	Definições e Revisões de Cálculo	4
1.3	Soluções Locais: Condições de primeira ordem	7
1.4	Soluções Locais: Condições de segunda ordem	7
1.5	Existência de pontos ótimos	11
1.6	Condições para soluções globais	13
1.7	Funções quadráticas	13
2	Otimização Convexa	16
2.1	Convexidade	17
2.1.1	Caracterização de convexidade de primeira ordem	20
2.1.2	Caracterizações de convexidade de segunda ordem	22
2.1.3	Convexidade forte	23
2.2	Otimização sobre conjuntos convexos	24
2.2.1	Condição de primeira ordem: Caso geral	24
2.2.2	Condições de primeira ordem: Caso convexo	27
3	Otimização com restrições lineares	28
3.1	Condições KKT	29
3.2	Condições KKT: Problema convexo	29
3.3	Condições KKT com restrições lineares de igualdade	30
4	Otimização com restrições genéricas	32
4.1	Lagrangeano	33
4.2	As generalizações do KKT	33
4.3	Caso Convexo	34
5	Algoritmos de Otimização	37
5.1	Método Gradiente	38

Otimização Irrestrita

Nota: Esse capítulo será cheio de definições e teoremas um atrás do outro, já que ele é mais uma revisão de coisas que já vimos em cursos passados. Todas as definições, teoremas e provas aqui escritas podem ser encontradas nas anotações do professor (Phillip Thompson). O intuito desse documento é esclarecer alguns conceitos que podem parecer confusos e dar intuições para vários conceitos bastante abstratos

1.1 Introdução

Otimização é um ramo da matemática preocupada em resolver problemas em que você possui várias opções de escolha, de forma que cada uma tem o custo associado, e queremos escolher a escolha com menor custo possível, ou seja, queremos resolver:

$$\min_{x \in C} f(x) \quad (1)$$

Com $f : C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ sendo a **função objeto** e C sendo o **conjunto viável**.

1.2 Definições e Revisões de Cálculo

Nesse capítulo, vamos rever alguns conceitos de cálculo e introduzir a otimização irrestrita, onde queremos trabalhar em uma função f sem nenhuma restrição

Definição 1.2.1 (Ponto de Mínimo): Seja $f : C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$

- $x^* \in C$ é um **ponto de mínimo global** de f em $C \Leftrightarrow$

$$\forall x \in C, \quad f(x^*) \leq f(x) \quad (2)$$

- $x^* \in C$ é um **ponto de mínimo global estrito** de f em $C \Leftrightarrow$

$$\forall x \in C, \quad f(x^*) < f(x) \quad (3)$$

- $x^* \in C$ é um **ponto de mínimo local** de f em $C \Leftrightarrow$

$$\exists r > 0 \wedge \forall x \in C \cap B(x^*, r), \quad f(x^*) \leq f(x) \quad (4)$$

- $x^* \in C$ é um **ponto de mínimo local estrito** de f em $C \Leftrightarrow$

$$\exists r > 0 \wedge \forall x \in C \cap B(x^*, r) \setminus \{x^*\}, \quad f(x^*) < f(x) \quad (5)$$

Definição 1.2.2 (Bola): Uma bola $B \subset \mathbb{R}^n$ de raio $r > 0 \in \mathbb{R}$ e centro $p \in \mathbb{R}^n$ é o conjunto:

$$B(p, r) = \{x \in \mathbb{R}^n / \|x - p\| \leq r\} \quad (6)$$

Ou seja, qual é a diferença dos dois? Pontos de mínimo globais são menores que todo e qualquer outro ponto no domínio C da função, enquanto os pontos locais são os menores em uma determinada vizinhança, a partir do ponto de mínimo local em questão, qualquer direção que eu tomar eu vou começar a subir o valor de f , mesmo que existam pontos em outros locais do domínio que sejam menores que o ponto de mínimo local que eu estava analisando

Agora vamos lembrar algumas coisas que vimos em cálculo (Alguns teoremas que são apenas revisão não serão demonstrados)

Definição 1.2.3 (Derivada direcional): Dada $f : \mathbb{R}^n \rightarrow \mathbb{R}$ e $d \neq 0 \in \mathbb{R}^n$ e $\|d\| = 1$. Se

$$\exists \lim_{t \rightarrow 0^+} \frac{f(x + td) - f(x)}{t} \quad (7)$$

Isso é chamado de derivada direcional de f na direção d ($\frac{df}{dd}$)

Definição 1.2.4 (Gradiente): Dada $f : \mathbb{R}^n \rightarrow \mathbb{R}$ e $\exists \frac{\partial f}{\partial x_i}, i = 1, \dots, n$, o vetor gradiente de f é definido como:

$$\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \quad (8)$$

Definição 1.2.5 (Continuamente Diferenciável): Uma função $f : \mathbb{R}^n \rightarrow \mathbb{R}$ é continuamente diferenciável se:

$$\forall x \in \mathbb{R}^n, \exists \frac{\partial f}{\partial x_i}(x) \quad (9)$$

e $\frac{\partial f}{\partial x_i}$ são contínuas $\forall i$

Teorema 1.2.1 (Aproximação de Primeira Ordem): Quando f é continuamente diferenciável, em uma vizinhança de um ponto x podemos mostrar que:

$$\forall d \in \mathbb{R}^n \text{ com } \|d\| = 1 \quad \frac{df}{dd} = \nabla f(x)^T d \quad (10)$$

e, além disso, temos:

$$\forall y \in \mathbb{R}^n \text{ na vizinhança} \quad f(y) = f(x) + \nabla f(x)^T (y - x) + o(\|y - x\|) \quad (11)$$

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

Apenas para relembrar, esse teorema está nos dando uma forma de aproximar uma função:

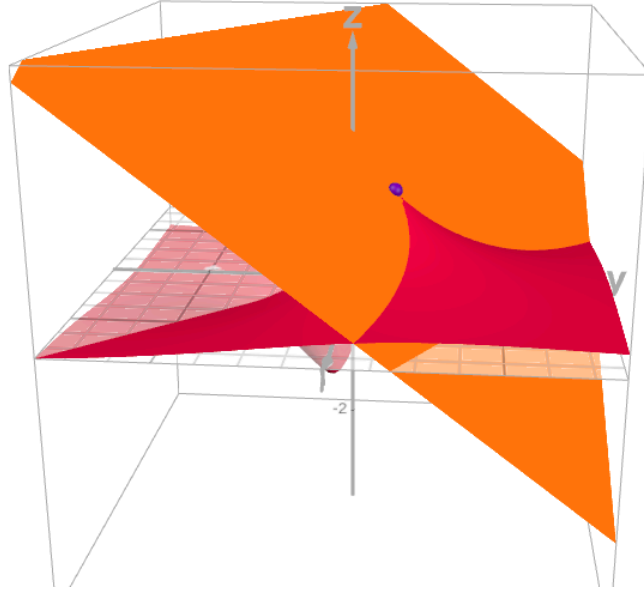


Figura 1: Função $f(x, y) = \frac{x+y}{x^2+y^2+1/5}$

Perceba que, próximo do ponto, a distância entre os pontos da curva e os do plano não são tão grandes, por isso que definimos a aproximação linear como mostrado anteriormente

Definição 1.2.6 (Funções duas vezes continuamente diferenciáveis): Podemos também expressar uma definição similar para uma função $f : C \rightarrow \mathbb{R}$ definida num conjunto $C \subset \mathbb{R}^n$. Dizemos que $f : C \rightarrow \mathbb{R}$ é duas vezes continuamente diferenciável em C se existe $U \supset C$ conjunto aberto tal que existem todas derivadas parciais de primeira e segunda ordem em todo ponto $x \in U$ e, além disso, as funções $\frac{\partial^2 f}{\partial x_i \partial x_j} : U \rightarrow \mathbb{R}$ são contínuas

Definição 1.2.7 (Matriz Hessiana): Seja $f : U \rightarrow \mathbb{R}$ com $U \subset \mathbb{R}^n$ e duas vezes continuamente diferenciável, a matriz hessiana de f no ponto $x \in U$ é definida como:

$$\nabla^2 f(x) := \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \ddots & \vdots \\ \vdots & & & \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix} \quad (12)$$

Perceba que $\nabla^2 f(x)$ é simétrica

Teorema 1.2.2 (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:

$$\forall y \in B(x, r) \exists \xi \in [x, y] \text{ tal que} \quad f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} (y - x)^T \nabla^2 f(\xi) (y - x) \quad (13)$$

Teorema 1.2.3 (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:

$$\forall y \in B(x, r) \text{ vale} \quad f(y) = f(x) + \nabla f(x)^T(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x)(y - x) + o(\|y - x\|^2) \quad (14)$$

1.3 Soluções Locais: Condições de primeira ordem

Agora podemos começar a brincadeira. Quando falamos de condições de primeira ordem, estamos nos referindo a condições relacionadas a derivadas de primeiro grau, ou seja, funções que são continuamente diferenciáveis. Antes eu comentei que estávamos interessados em minimizar funções num conjunto C , porém, vamos primeiro ver sobre otimização **irrestrita**, ou seja, problemas do tipo:

$$\min_{x \in \mathbb{R}^n} f(x) \quad (15)$$

Lembram que o vetor gradiente indica a direção que minha função tá crescendo? Quando estamos procurando um mínimo local, faz sentido dizer que a função cresça pra todos os lados, correto? Então faz sentido dizer que isso vai me dar um vetor gradiente 0 (Apenas uma intuição)

Teorema 1.3.1 (Condições de primeira ordem): Seja $f : U \rightarrow \mathbb{R}$ uma função definida no conjunto aberto $U \subset \mathbb{R}^n$. Se $x^* \in U$ é um mínimo local de f e todas as derivadas parciais de f existem, então

$$\nabla f(x^*) = 0 \quad (16)$$

Demonstração: Seja $i \in [n]$ e defina a função $g(t) = f(x^* + te_i)$. Temos que g é diferenciável em 0 e $g'(0) = \frac{\partial f}{\partial x_i}(x^*)$. Sendo x^* um ponto ótimo local de f , segue que 0 é um ponto ótimo local de g ; portanto $0 = g'(0) = \frac{\partial f}{\partial x_i}(x^*)$. O argumento vale para todo $i \in [n]$, implicando que $\nabla f(x^*) = 0$ \square

Esse teorema não vale na volta, já que, como vimos antes em cálculo, pontos de máximo e de sela também possuem essa característica, isso nos leva a criar a definição:

Definição 1.3.1 (Ponto estacionário): Seja $f : U \rightarrow \mathbb{R}$ uma função definida no conjunto aberto $U \subset \mathbb{R}^n$ e todas as derivadas parciais de f existem, então chamamos $x^* \in U$ de ponto estacionário de f em U se

$$\nabla f(x^*) = 0 \quad (17)$$

1.4 Soluções Locais: Condições de segunda ordem

Nas anotações, o professor generaliza o conceito de que, se x é estacionário e $f''(x) > 0$ então x é mínimo local.

Definição 1.4.1 (Positividade e Negatividade de uma matriz): Seja A uma matriz simétrica:

- Dizemos que A é positiva semidefinida, denotando-s por $A \succeq 0 \Leftrightarrow \forall x \in \mathbb{R}^n, x^T A x \geq 0$
- Dizemos que A é positiva definida, denotando-s por $A \succ 0 \Leftrightarrow \forall x \in \mathbb{R}^n, x^T A x > 0$
- Dizemos que A é negativa semidefinida, denotando-s por $A \preceq 0 \Leftrightarrow \forall x \in \mathbb{R}^n, x^T A x \leq 0$
- Dizemos que A é negativa definida, denotando-s por $A \prec 0 \Leftrightarrow \forall x \in \mathbb{R}^n, x^T A x < 0$
- Dizemos que A é indefinida, quando $\exists x, y \in \mathbb{R}^n$ tal que $x^T A x > 0$ e $y^T A y < 0$

Nas anotações do professor ele traz alguns conceitos que vimos em álgebra linear, mas eu não vou os abordar aqui.

Teorema 1.4.1 (Condições necessárias de segunda ordem): Seja $f : U \rightarrow \mathbb{R}$ com $U \subset \mathbb{R}^n$ e suponha que f é duas vezes continuamente diferenciável sobre U e seja $x^* \in U$, então:

1. Se x^* é mínimo local, então $\nabla^2 f(x^*) \succeq 0$
2. Se x^* é máximo local, então $\nabla^2 f(x^*) \preceq 0$

Demonstração: Vamos apenas provar o item 1 já que a prova para o 2 é análoga (Basta aplicar a demonstração na função $-f$).

Sendo x^* um ponto de mínimo local, $\exists B(x^*, r) \subset U$ tal que:

$$\forall x \in B(x^*, r) \quad f(x) \geq f(x^*) \quad (18)$$

Seja $0 \neq d \in \mathbb{R}^n$. Para todo $0 < \alpha < \frac{r}{\|d\|}$, vamos definir:

$$x_\alpha^* := x^* + \alpha d \quad (19)$$

$$x_\alpha^* \in B(x^*, r) \Rightarrow f(x_\alpha^*) \geq f(x^*) \quad (20)$$

Pelo Teorema 1.2.2, $\exists \xi_\alpha \in [x^*, x_\alpha^*]$ tal que

$$f(x_\alpha^*) - f(x^*) = \nabla f(x^*)^T (x_\alpha^* - x^*) + \frac{1}{2} (x_\alpha^* - x^*)^T \nabla^2 f(\xi_\alpha) (x_\alpha^* - x^*) \quad (21)$$

Como x^* é estacionário, temos:

$$f(x_\alpha^*) - f(x^*) = \frac{\alpha^2}{2} d^T \nabla^2 f(\xi_\alpha) d \quad (22)$$

Combinando as equações (18) e (22), temos que, para todo $\alpha \in (0, \frac{r}{\|d\|})$:

$$d^T \nabla^2 f(\xi_\alpha) d \geq 0 \quad (23)$$

Usando de que $\xi_\alpha \rightarrow x^*$ quando $\alpha \rightarrow 0^+$ e por continuidade da Hessiana, segue:

$$d^T \nabla^2 f(x^*) d \geq 0 \quad (24)$$

Isso é válido pois eu assumi um d genérico □

Perceba que essa condição é necessária, mas não é suficiente. Por exemplo, a função $f(x) = x^3$ é tal que $f'(0) = 0$, $f''(0) = 0$, porém, não é um ponto de máximo nem de mínimo.

Teorema 1.4.2 (Condições suficientes de segunda ordem): Seja $f : U \rightarrow \mathbb{R}$ com $U \subset \mathbb{R}^n$ e suponha que f é duas vezes continuamente diferenciável sobre U e seja $x^* \in U$ um ponto estacionário de f em U , então:

- Se $\nabla^2 f(x^*) \succ 0$ então x^* é um ponto de mínimo local estrito
- Se $\nabla^2 f(x^*) \prec 0$ então x^* é um ponto de máximo local estrito

Demonstração: Provaremos apenas o primeiro item. O segundo segue do primeiro aplicado em $-f$

Seja $x^* \in U$ um ponto estacionário de f em U tal que $\nabla^2 f(x^*) \succ 0$. Como a Hessiana é contínua, segue que $\exists B(x^*, r) \subset U$ tal que $\nabla^2 f(x^*) \succ 0 \forall x \in B(x^*, r)$. Pelo Teorema 1.2.2, segue que $\forall x \in B(x^*, r) \exists \xi \in [x^*, x] \subset B(x^*, r)$ tal que:

$$f(x) - f(x^*) = \nabla f(x^*)^T (x - x^*) + \frac{1}{2}(x - x^*)^T \nabla^2 f(\xi)(x - x^*) \quad (25)$$

Como x^* é estacionário, $\nabla f(x^*) = 0$. Segue também que $\nabla^2 f(\xi) \succ 0 \forall x \in B(x^*, r)$. Isso significa que

$$\forall x \neq x^*, \quad f(x) > f(x^*) \quad (26)$$

Ou seja, x^* é mínimo local estrito □

Para clarear um pouco sobre a demonstração, os passos mais confusos pode ser a conclusão final. Principalmente essa conclusão $\nabla^2 f(\xi) \succ 0$. Vamos tentar abstrair isso com $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Pega um ponto de mínimo estrito local, e faz uma bola em volta dele, todo ponto dentro daquele lugar vai ter hessiana positiva por conta da continuidade da Hessiana. Como assim? Imagina que a Hessiana é uma função $\mathbb{R} \rightarrow \mathbb{R}$, como ela é uma função contínua, não faz sentido eu mudar a entrada da função e ela bruscamente trocar de positivo pra negativo, certo? Claro que em um certo ponto, ela passa pelo 0 e o sinal troca, mas eu consigo aumentar minha bola até um pouquinho antes disso acontecer

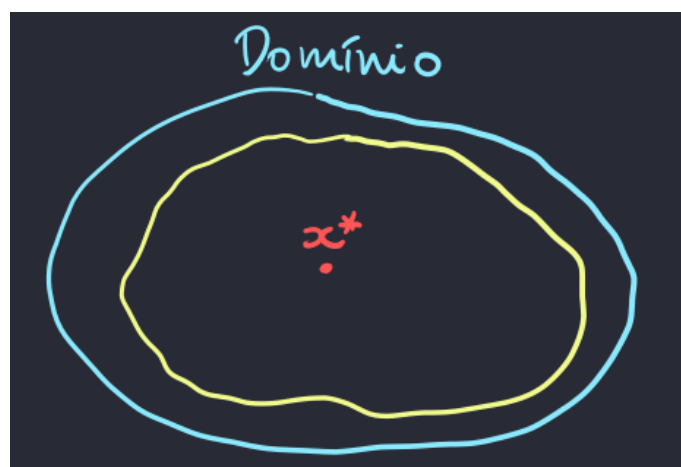


Figura 2: Desenho de domínio qualquer de uma função f

A partir dessa linha amarela, os pontos vão ter hessiana negativa e, em cima dela, eles tem hessiana igual a 0, ou seja, então eu consigo criar uma bola $B(x^*, r)$ de forma que ela não ultrapasse a linha amarela

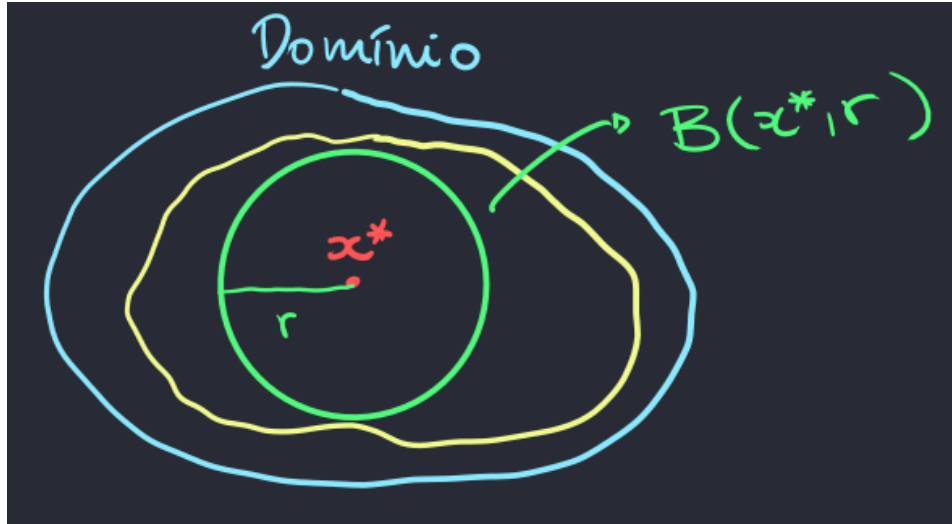


Figura 3: Desenho de domínio qualquer de uma função f com uma bola B

Ou seja, eu sei que todos os pontos dentro dessa bola tem Hessiana positiva. Depois disso, eu apenas utilizo do Teorema 1.2.2 para chegar na desigualdade $f(x) > f(x^*)$

Um teorema parecido pode ser usado para pontos que tem gradiente 0, mas que não são nem máximo nem mínimo (Como vimos em $f(x) = x^3$)

Definição 1.4.2 (Ponto de Sela): Seja $f : U \rightarrow \mathbb{R}$ definida num conjunto aberto $U \subset \mathbb{R}^n$. Suponha que f é duas vezes continuamente diferenciável. $x^* \in U$ é ponto de sela de f em U se ele é um ponto estacionário, mas não é nem ponto de máximo nem ponto de mínimo

Teorema 1.4.3 (Condições suficientes para pontos de sela): Seja $f : U \rightarrow \mathbb{R}$ com $U \subset \mathbb{R}^n$ e suponha que f é duas vezes continuamente diferenciável sobre U e seja $x^* \in U$ um ponto estacionário de f em U , se $\nabla^2 f(x^*)$ é indefinida, então x^* é ponto de sela

Demonstração: Seja $\nabla^2 f(x^*)$ é indefinida. Portanto, $\nabla^2 f(x^*)$ possui auto-valor positivo λ_1 associado ao auto-vetor v_1 com norma $\|v_1\| = 1$. Sendo U aberto, existe $r > 0$ tal que $x^* + \alpha v_1 \in U$ para todo $\alpha \in (0, r)$. Pelo Teorema 1.2.3 e usando que $\nabla f(x^*) = 0$, sabemos que existe uma função $o : \mathbb{R}_+ \rightarrow \mathbb{R}$ satisfazendo:

$$\lim_{t \rightarrow 0^+} \frac{o(t)}{t} = 0 \quad (27)$$

tal que para todo $\alpha \in (0, r)$:

$$\begin{aligned} f(x^* + \alpha v_1) &= f(x^*) + \frac{\alpha^2}{2} v_1^T \nabla^2 f(x^*) v_1 + o(\alpha^2 \|v_1\|^2) \\ &= f(x^*) + \frac{\lambda_1 \alpha^2}{2} \|v_1\|^2 + o(\alpha^2 \|v_1\|^2) \\ &= f(x^*) + \frac{\lambda_1 \alpha^2}{2} + o(\alpha^2) \end{aligned} \quad (28)$$

Segue da equação (27) que $\exists \varepsilon_1 \in (0, r)$ tal que:

$$\forall \alpha \in (0, \varepsilon_1), \quad g(\alpha^2) > -\frac{\lambda_1 \alpha^2}{2} \quad (29)$$

Portanto,

$$\forall \alpha \in (0, \varepsilon_1), \quad f(x^* + \alpha v_1) > f(x^*) \quad (30)$$

Ou seja, x^* não pode ser máximo local sobre U . Um argumento análogo dizendo que $\exists \lambda_2 < 0$ sendo λ_2 um autovalor da hessiana pode ser usado para mostrar que x^* também não pode ser mínimo local \square

Essa prova parece complicada, então vou dar uma noção mais intuitiva. Vimos em álgebra linear que uma matriz é positiva definida se, e somente se, todos os seus autovalores são maiores que 0 (O mesmo para matrizes negativas definidas), e que se elas possuem um autovalor positivo e outro negativo, então ela é indefinida. Mas o que isso me diz intuitivamente? Lembra que, se uma matriz tem multiplicidade algébrica igual a multiplicidade geométrica em todos os autovalores, então a gente pode dividir ela como:

$$\nabla^2 f(x) = Q^T \Lambda Q \quad (31)$$

Q é ortogonal pois $\nabla^2 f(x)$ é simétrica (Teorema Espectral). Mas o que isso significa? De uma maneira intuitiva, isso significa que os autovetores indicam direções ortogonais e o autovalor indica se a hessiana está crescendo ou diminuindo **naquela direção**, então se ela é indefinida em um ponto de sela, quer dizer que eu tenho direções que a hessiana tanto cresce como diminui, como ela cresce e diminui em direções diferentes partindo do mesmo ponto, ele não é nem máximo, nem mínimo

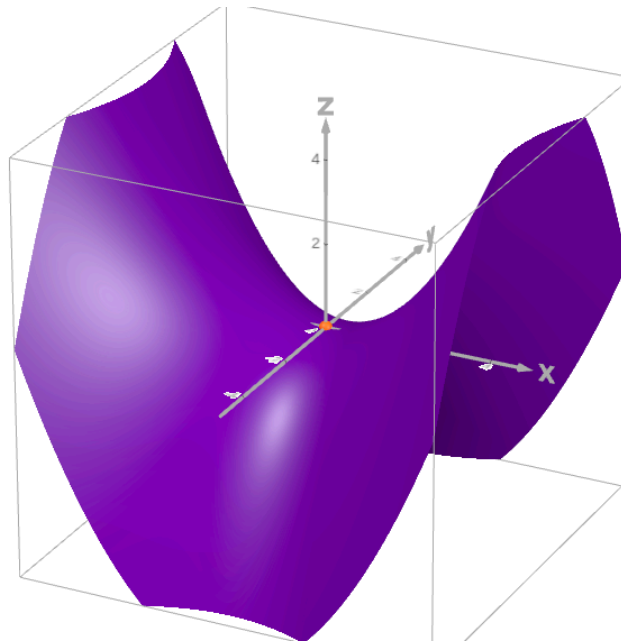


Figura 4: Função $f(x, y) = ax^2 + by^2$. Ponto laranja é ponto de sela (Ponto $(0,0,0)$)

1.5 Existência de pontos ótimos

Até agora estávamos assumindo que pontos ótimos existiam, mas e se eles não existem?

Definição 1.5.1 (Conjunto fechado): Um conjunto C é fechado se seu complementar C^c é aberto

Definição 1.5.2 (Conjunto limitado): Um conjunto C é limitado se $\exists r > 0$ tal que $C \subset B(0, r)$

Definição 1.5.3 (Conjunto compacto): Um conjunto C é compacto se é fechado e limitado

Teorema 1.5.1 (Weierstrass): Seja $C \subset \mathbb{R}^n$ um conjunto compacto e $f : C \rightarrow \mathbb{R}$, então f possui um ponto de mínimo global e de máximo global em C

Quando o conjunto não é compacto, o teorema de Weierstrass não garante a existência, então podemos usar essa outra definição:

Definição 1.5.4 (Coercividade): Seja $f : \mathbb{R}^n \rightarrow \mathbb{R}$. A função é dita coerciva se:

$$\lim_{\|x\| \rightarrow \infty} f(x) = \infty \quad (32)$$

Ou seja, todo e qualquer vetor que eu pegar e aumentar seu tamanho, a função aumenta junto, formando o que parece uma grande bacia, onde você coloca água e ela nunca vaza

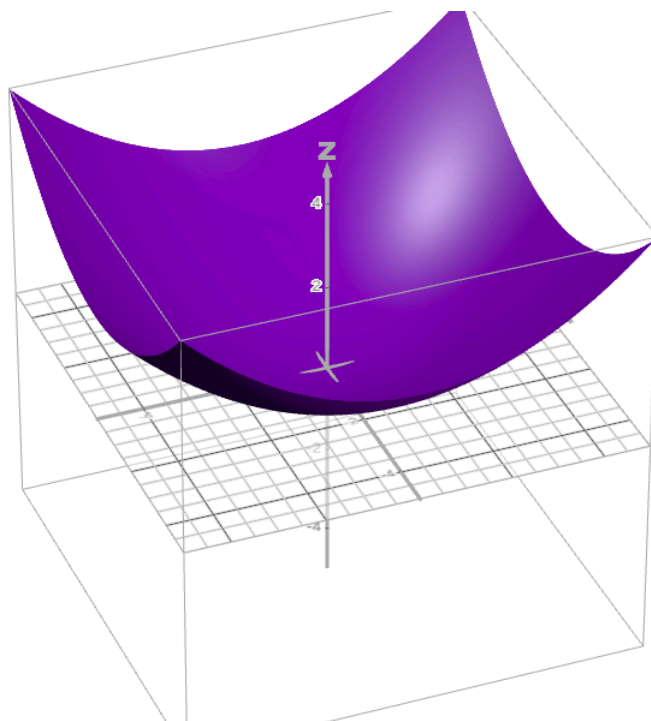


Figura 5: Exemplo de função coerciva $f(x, y) = 0.1x^2 + 0.1y^2$

Teorema 1.5.2 (Existência de soluções: Coercividade): Seja $f : \mathbb{R}^n \rightarrow \mathbb{R}$ uma função contínua e coerciva e $C \subset \mathbb{R}^n$ um conjunto fechado não-vazio. Então f tem um mínimo global em C

Demonstração: Seja $x_0 \in C$ um ponto arbitrário. Como f é coerciva, segue que existe $M > 0$ tal que

$$f(x) > f(x_0) \text{ para todo } x \text{ tal que } \|x\| > M \quad (33)$$

Temos que x^* é um ponto de mínimo global de f sobre C . Portanto $f(x^*) \geq f(x_0)$. Segue da afirmação em display que o conjunto de mínimos globais de f sobre C é exatamente o conjunto de mínimos globais de f sobre $C \cap B(0, M)$. O conjunto $C \cap B(0, M)$ é fechado e limitado, portanto compacto. Segue do Teorema de Weierstrass que f possui ponto de mínimo global sobre $C \cap B(0, M)$, e portanto, sobre C também \square

1.6 Condições para soluções globais

Teorema 1.6.1: Seja $f : \mathbb{R}^n \rightarrow \mathbb{R}$ duas vezes continuamente diferenciável. Suponha que:

$$\nabla^2 f(x) \succeq 0, \quad \forall x \in \mathbb{R}^n \quad (34)$$

Então, em todo ponto estacionário de f , esse ponto é um mínimo global

Demonstração: Pelo Teorema 1.2.2, seja $x^* \in \mathbb{R}^n$ um ponto estacionário em f e $\forall x \in \mathbb{R}^n$:

$$f(x) - f(x^*) = \frac{1}{2}(x - x^*)^T \nabla^2 f(\xi)(x - x^*) \quad (35)$$

Porém, vale que $\forall x, \nabla^2 f(\xi) \succeq 0$. Temos então que:

$$\forall x \in \mathbb{R}^n, f(x) \geq f(x^*) \quad (36)$$

Logo, x^* é ponto de mínimo global em f \square

1.7 Funções quadráticas

Um conjunto interessante de funções com algumas propriedades convenientes são as funções quadráticas

Definição 1.7.1 (Função quadrática): Uma função é quadrática quando $\exists A \in \mathbb{R}^{n \times n}$ simétrica, $b \in \mathbb{R}^n, c \in \mathbb{R}$ tal que a função $f : \mathbb{R}^n \rightarrow \mathbb{R}$ pode ser expressa como:

$$f(x) = x^T A x + 2b^T x + c \quad (37)$$

Teorema 1.7.1 (Derivadas de uma quadrática): Seja f uma função quadrática como na Definição 1.7.1, temos que:

$$\begin{aligned} \nabla f(x) &= 2(Ax + b) \\ \nabla^2 f(x) &= 2A \end{aligned} \quad (38)$$

Demonstração: Sabemos que $f(x) = x^T A x + 2b^T x + c$. Vamos definir que x_i é a i -ésima entrada de x . Vamos primeiro calcular uma derivada parcial genérica de f . Como a derivada é uma operação linear, eu vou ver cada componente separadamente.

$$x^T A x = (x_1 \dots x_n) \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = (x_1 \dots x_n) \begin{pmatrix} \sum_{k=1}^n a_{1k} x_k \\ \vdots \\ \sum_{k=1}^n a_{nk} x_k \end{pmatrix} \quad (39)$$

Para facilitar nossa vida, vamos definir

$$\alpha_j = \sum_{k=1}^n a_{jk} x_k \quad (40)$$

. Então:

$$f(x) = \alpha_1 x_1 + \dots + \alpha_n x_n + 2(b_1 x_1 + \dots + b_n x_n) + c \quad (41)$$

Agora podemos tirar a derivada de $f(x)$ em x_j , mas antes, perceba que:

$$\frac{\partial \alpha_i}{\partial x_j} = a_{ij} \quad (42)$$

Agora sim:

$$\begin{aligned} \frac{\partial f}{\partial x_j} &= x_1 \frac{\partial \alpha_1}{\partial x_j} + \dots + \frac{\partial}{\partial x_j}(\alpha_j x_j) + \dots + x_n \frac{\partial \alpha_n}{\partial x_j} + 2b_j \\ \frac{\partial f}{\partial x_j} &= x_1 a_{1j} + \dots + \frac{\partial \alpha_j}{\partial x_j} x_j + \alpha_j + \dots + x_n a_{nj} + 2b_j \\ \frac{\partial f}{\partial x_j} &= \sum_{k=1}^n a_{jk} x_k + \sum_{k=1}^n a_{kj} x_k + 2b_j \end{aligned} \quad (43)$$

Como A é simétrica, podemos reescrever isso como:

$$\frac{\partial f}{\partial x_j} = 2 \left(\sum_{k=1}^n a_{kj} x_k + b_j \right) \quad (44)$$

Ou seja, o gradiente da função é:

$$\nabla f(x) = 2(Ax + b) \quad (45)$$

E para a hessiana é bem mais fácil, dado o item anterior, basta que tiremos a derivada novamente para x_i :

$$\frac{\partial^2 f}{\partial x_j \partial x_i} = 2a_{ij} \quad (46)$$

Ou seja:

$$\nabla^2 f(x) = 2A \quad (47)$$

□

Teorema 1.7.2 (Pontos estacionários e ótimos de função quadrática): Seja uma função f definida na Definição 1.7.1, então:

1. x é ponto estacionário $\Leftrightarrow Ax = -b$.
2. Suponha que $A \succeq 0$. Então x é ponto de mínimo global $\Leftrightarrow Ax = -b$.
3. Suponha que $A \succ 0$. Então $x = -A^{-1}b$ é ponto de mínimo global estrito.

Demonstração:

1. Segue imediatamente da fórmula do gradiente.

2. Suponha que $A \succeq 0$. Da fórmula da Hessiana, segue que $\nabla^2 f(x) \succeq 0 \forall x \in \mathbb{R}^n$. O resultado segue então do Teorema 1.6.1 e item 1.
3. Suponha que $A \succ 0$. Então $x = -A^{-1}b$ é a única solução de $Ax = -b$. Segue do item (ii) que $x = -A^{-1}b$ é o único ponto de mínimo global de f e, portanto, mínimo global estrito.

□

Teorema 1.7.3 (Coercividade de funções quadráticas): Seja função f definida como na Definição 1.7.1. Então f é coerciva $\Leftrightarrow A \succ 0$.

Demonstração: Precisamos do seguinte lema: Seja $A \in \mathbb{R}^{n \times n}$ simétrica, então $\forall x \neq 0 \in \mathbb{R}^{n \times n}$

$$\lambda_{\min}(A) \leq \frac{x^T A x}{\|x\|^2} \leq \lambda_{\max}(A) \quad (48)$$

(Pode-se demonstrar pelo teorema espectral)

Agora podemos começar a prova:

(\Leftarrow) Suponha que $A \succ 0$. Denote $\alpha := \lambda_{\min}(A)$. Pelo lema acima e Cauchy-Schwarz, segue que, para todo $x \in \mathbb{R}^n$,

$$f(x) = x^T A x + 2b^T x + c \geq \alpha \|x\|^2 - 2\|b\|\|x\| + c \quad (49)$$

Segue que $f(x) \rightarrow \infty$ quando $\|x\| \rightarrow \infty$; isto é, f é coerciva.

(\Rightarrow) Suponha que f é coerciva. Suponha que A tenha auto-valores negativos. Portanto, existem $v \neq 0$ e $\lambda < 0$ tais que $Av = \lambda v$. Portanto, para todo $\alpha \in \mathbb{R}$,

$$f(\alpha v) = \lambda \|\alpha v\|^2 + 2(b^T v)\alpha + c \rightarrow \infty \text{ quando } \alpha \rightarrow \infty \quad (50)$$

Isto contradiz a hipótese de coercividade. Portanto, A possui todos auto-valores não-negativos. Provaremos agora que 0 não é auto-valor de A , provando que $A \succ 0$. Assuma que exista $v \neq 0$ tal que $Av = 0$. Então, para todo $\alpha \in \mathbb{R}$,

$$f(\alpha v) = 2(b^T v)\alpha + c. \quad (51)$$

Temos que:

$$f(\alpha v) \rightarrow \begin{cases} c & \text{quando } \alpha \rightarrow \infty \text{ se } b^T v = 0 \\ -\infty & \text{quando } \alpha \rightarrow -\infty \text{ se } b^T v > 0 \\ \infty & \text{quando } \alpha \rightarrow \infty \text{ se } b^T v < 0 \end{cases} \quad (52)$$

Em qualquer caso a coerção é violada, portanto, 0 não pode ser autovalor de A

□

Otimização Convexa

Agora vamos focar os nossos esforços em resolver um conjunto específico de problemas de otimização, os do tipo:

$$\min_{x \in C} f(x) \quad C \subseteq \mathbb{R}^n \text{ convexo} \quad (53)$$

Ou seja, agora começaremos a aplicar as famosas restrições nos problemas que vamos abordar

2.1 Convexidade

Definição 2.1.1 (Conjunto convexo): Um conjunto $C \subseteq \mathbb{R}^n$ é dito convexo se

$$\forall x, y \in C \wedge \forall \lambda \in (0, 1) \text{ vale } \lambda x + (1 - \lambda)y \in C \quad (54)$$

Ou seja, se eu pego dois pontos dentro do conjunto C e fizer uma reta que interliga eles, todos os pontos nessa reta devem estar dentro de C

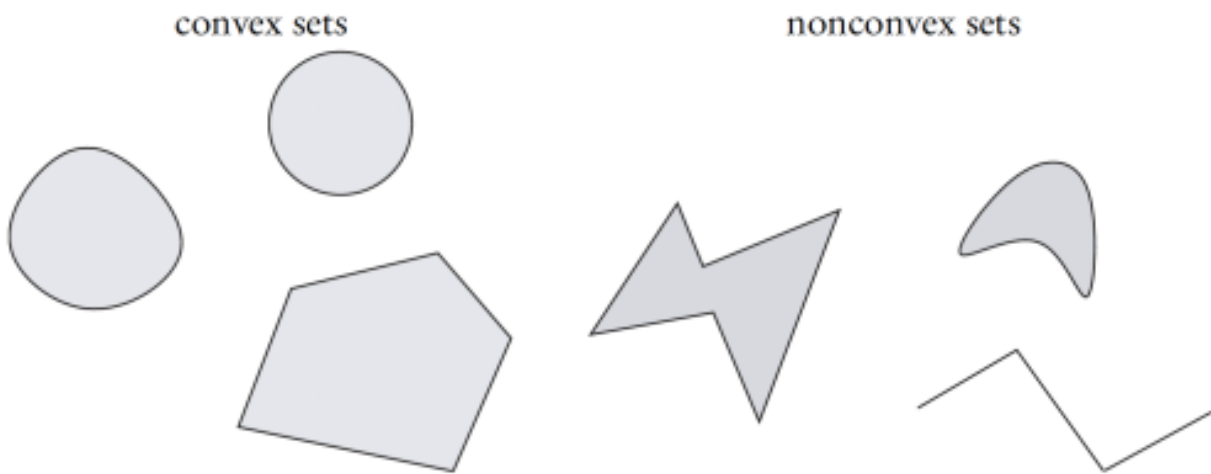


Figura 6: Exemplo de figuras convexas e não-convexas retirado das anotações do professor

Porém, outra definição muito importante são as de **funções convexas**

Definição 2.1.2 (Funções convexas): Uma função $f : C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ com C convexo é dita convexa se:

$$\forall x, y \in C \wedge \forall \lambda \in (0, 1) \text{ vale } f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (55)$$

Definição 2.1.3 (Funções estritamente convexas): Uma função $f : C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$ com C convexo é dita estritamente convexa se:

$$\forall x, y \in C \wedge \forall \lambda \in (0, 1) \text{ vale } f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y) \quad (56)$$

Mas o que isso quer dizer? Quer dizer que eu vou pegar o segmento entre meus pontos x e y e vou aplicar a função neles, depois eu vou pegar o segmento de reta entre $f(x)$ e $f(y)$ e comparar. Todos os pontos nesse segmento de reta tem que estar acima dos pontos da curva que eu fiz antes

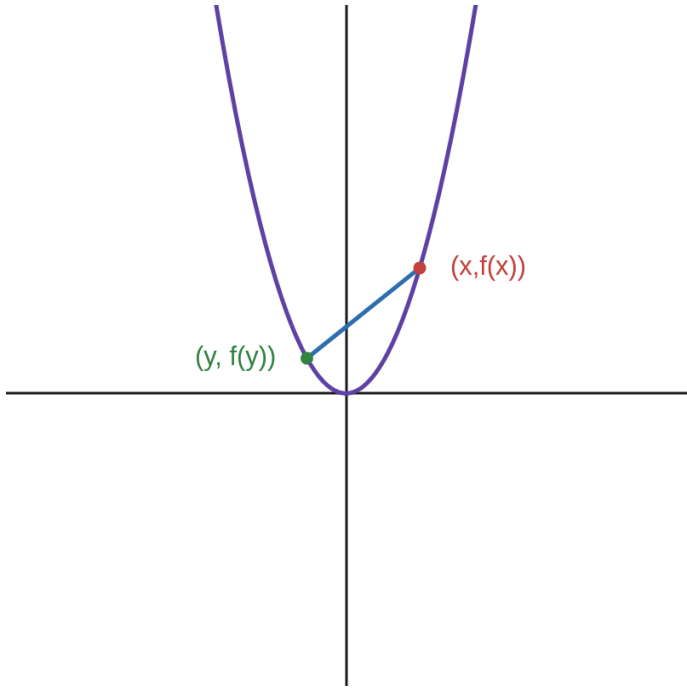


Figura 7: Função convexa $f(x) = x^2$

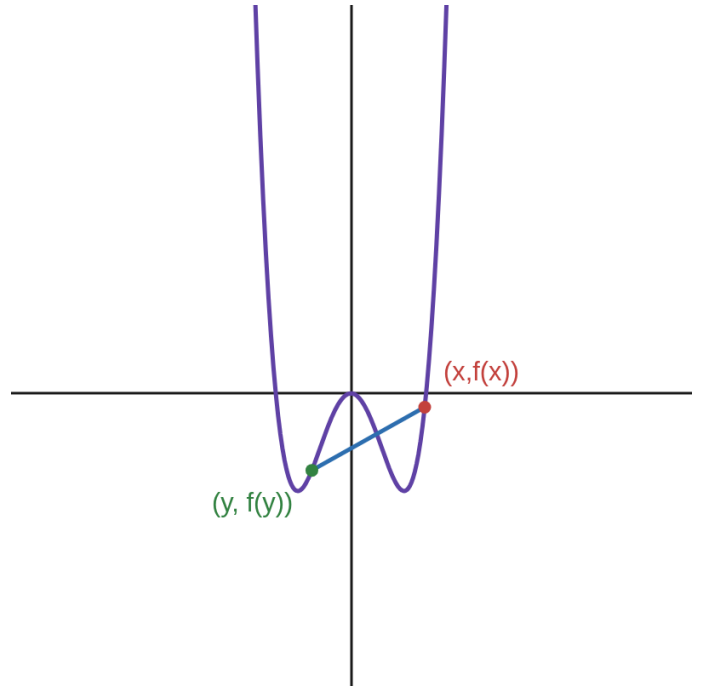


Figura 8: Função não-convexa $f(x) = x^4 - 3x^2$

Antes de continuar, vamos definir um conjunto simplex, que utilizaremos bastante daqui pra frente:

Definição 2.1.4 (Conjunto simplex): O conjunto simplex Δ_k é definido como:

$$\Delta_k := \left\{ \lambda \in \mathbb{R}_+^k / \sum_{i=1}^k \lambda_i = 1 \right\} \quad (57)$$

Existe um teorema que mostra que isso vale não só para a combinação de dois pontos, mas para a combinação de quaisquer n pontos

Teorema 2.1.1 (Teorema de Jenssen): Seja $f : C \subseteq \mathbb{R}^n \rightarrow \mathbb{R}$, com C convexo, uma função convexa. Então dados quais quer coleção $\{x_i\}_{i=1}^k \subset C$ de pontos de C e qualquer $\lambda \in \Delta_k$:

$$f\left(\sum_{i=1}^k \lambda_i x_i\right) \leq \sum_{i=1}^k \lambda_i f(x_i) \quad (58)$$

Demonstração: Faremos por indução. O caso base $k = 1$ é bem óbvio. Agora vamos supor que vale para k . Sejam $\{x_i\}_{i=1}^{k+1} \subset C$ e $\lambda \in \Delta_{k+1}$. Para facilitar, definamos:

$$z := \sum_{i=1}^{k+1} \lambda_i x_i \quad (59)$$

Se $\lambda_{k+1} = 1$, então $\sum_{i=1}^k \lambda_i = 0$; Como $\lambda_i \geq 0 \forall i \in [k]$, tem-se que $\lambda_i = 0$. Nesse caso, $x_{k+1} = z$ e a desigualdade é imediata. Se $\lambda_{k+1} < 1$. Nesse caso,

$$z = \lambda_{k+1} x_{k+1} + \underbrace{\left(\sum_{i=1}^k \lambda_i x_i \right)}_v = \lambda_{k+1} x_{k+1} + (1 - \lambda_{k+1}) \sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} x_i \quad (60)$$

E é bem fácil de ver que

$$\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} = 1 \quad (61)$$

Como C é convexo e $\{x_i\}_{i=1}^k \subset C$, então temos que:

$$\frac{1}{1 - \lambda_{k+1}} \sum_{i=1}^k \lambda_i x_i \in C \quad (62)$$

Isso é um teorema que vou enunciar posteriormente e demonstrar também. Como $x_{k+1} \in C$, pela convexidade de f ,

$$\begin{aligned} f(z) &= f((1 - \lambda_{k+1})v + \lambda_{k+1}x_{k+1}) \leq (1 - \lambda_{k+1})f(v) + \lambda_{k+1}f(x_{k+1}) \\ &= (1 - \lambda_{k+1})f\left(\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} x_i\right) + \lambda_{k+1}f(x_{k+1}) \\ &\leq (1 - \lambda_{k+1})f\left(\sum_{i=1}^k \frac{\lambda_i}{1 - \lambda_{k+1}} f(x_i)\right) + \lambda_{k+1}f(x_{k+1}) \\ &= \sum_{i=1}^{k+1} f(x_i) \end{aligned} \quad (63)$$

□

Teorema 2.1.2: Seja $C \subset \mathbb{R}^n$ convexo, dados quaisquer coleção $\{x_i\}_{i=1}^k \subset C$ de pontos em C e qualquer $\lambda \in \Delta_k$, então:

$$\sum_{i=1}^k \lambda_i x_i \in C \quad (64)$$

Demonstração: Caso base: $k = 2$ é trivial. Vamos supor que vale para um determinado k , então:

$$\sum_{i=1}^k \mu_i x_i \in C \quad (65)$$

Para qualquer $\mu \in \Delta_k$. Já que esse ponto está em C , vamos pegar um novo vetor x_{k+1} ainda em C . Já que ambos os vetores estão em C e ele é convexo, vale:

$$\forall \alpha \in (0, 1), \alpha \sum_{i=1}^k \mu_i x_i + (1 - \alpha)x_{k+1} \in C \quad (66)$$

Porém, perceba que

$$\alpha \sum_{i=1}^k \mu_i + (1 - \alpha) = \alpha + 1 - \alpha = 1 \quad (67)$$

Ou seja, se eu denotar $\lambda \in \mathbb{R}^{k+1}$ de tal forma que $\lambda_i = \alpha \mu_i$ para $i \in [k]$ e $\lambda_{k+1} = 1 - \alpha$ eu obtenho uma coleção de números tal que $\lambda \in \Delta_k$ e uma coleção $\{x_i\}_{i=1}^{k+1}$ tal que:

$$\sum_{i=1}^{k+1} \lambda_i x_i \in C \quad (68)$$

□

2.1.1 Caracterização de convexidade de primeira ordem

Funções convexas podem ser não-diferenciáveis. Funções convexas diferenciáveis possuem uma caracterização importante: hiperplanos tangentes ao seu gráfico são sempre estimativas abaixo da função.

Teorema 2.1.1.1 (Desigualdade do gradiente): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ com C convexa e f continuamente diferenciável, então:

$$f \text{ convexa} \Leftrightarrow \forall x, y \in C, f(x) + \nabla f(x)^T(y - x) \leq f(y) \quad (69)$$

Demonstração: (\Rightarrow) Suponha que f seja convexa. Sejam $x, y \in C \wedge \lambda \in [0, 1]$. A desigualdade enunciada vale trivialmente se $x = y$. Iremos então assumir que $x \neq y$. Da convexidade de f ,

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (70)$$

implicando que

$$\frac{f(x + \lambda(y - x)) - f(x)}{\lambda} \leq f(y) - f(x) \quad (71)$$

Tomando $\lambda \rightarrow 0^+$, obtemos

$$f'(x; y - x) = \lim_{\lambda \rightarrow 0^+} \frac{f(x + \lambda(y - x)) - f(x)}{\lambda} \leq f(y) - f(x) \quad (72)$$

Como f é continuamente diferenciável, $f'(x; y - x) = \nabla f(x)^T(y - x)$ e a desigualdade segue.

(\Leftarrow) Assuma que a desigualdade vale. Sejam $x, y \in C$ e $\lambda \in (0, 1)$. Defina $z = \lambda x + (1 - \lambda)y$. Temos:

$$x - z = z - (1 - \lambda)y - z = (1 - \lambda)(\lambda)(y - z) \quad (73)$$

À seguir, usaremos a desigualdade nos pares (x, z) e (y, z) . Temos

$$f(z) + \nabla f(z)^T(x - z) \leq f(x) \quad (74)$$

$$f(z) + \nabla f(z)^T(y - z) \leq f(y) \quad (75)$$

Multiplicando-se a primeira desigualdade por $(\lambda)(1 - \lambda)$ e usando a igualdade na segunda desigualdade, obtemos

$$\frac{\lambda}{1 - \lambda} f(z) + \frac{\lambda}{1 - \lambda} \nabla f(z)^T(x - z) \leq \frac{\lambda}{1 - \lambda} f(x) \quad (76)$$

$$f(z) - \frac{\lambda}{1 - \lambda} \nabla f(z)^T(x - z) \leq f(y) \quad (77)$$

Somando-se as duas desigualdades acima obtemos

$$\frac{\lambda}{1-\lambda}f(z) + f(z) \leq \frac{\lambda}{1-\lambda}f(x) + f(y) \quad (78)$$

Isto é

$$f(z) \leq \lambda f(x) + (1-\lambda)f(y) \quad (79)$$

Segue que f é convexa

□

Teorema 2.1.1.2 (Desigualdade do gradiente estrito): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ com C convexa e f continuamente diferenciável, então:

$$f \text{ estritamente convexa} \Leftrightarrow \forall x, y \in C, f(x) + \nabla f(x)^T(y - x) < f(y) \quad (80)$$

Demonstração: Análogo ao Teorema 2.1.1.1

□

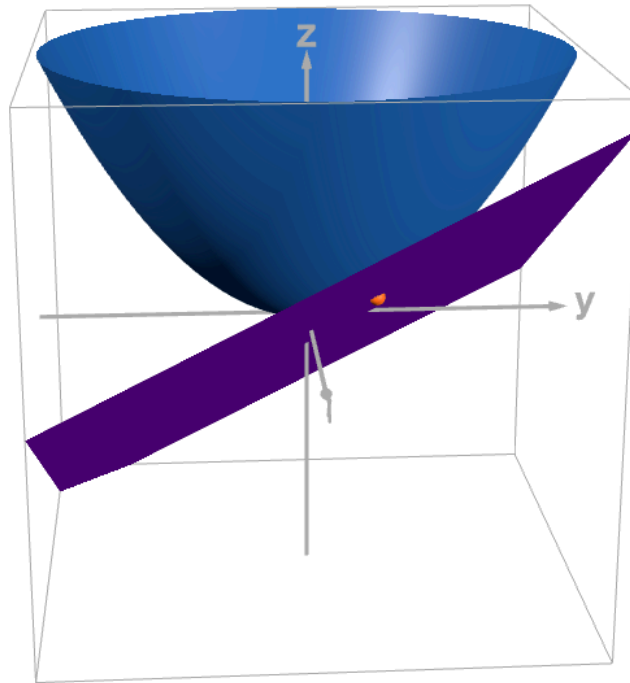


Figura 9: Função $f(x, y) = 1.3x^2 + 1.27y^2$ e um plano tangente à curva

A gente pode usar os teoremas anteriores pra caracterizar as funções quadráticas e quando elas são convexas

Teorema 2.1.1.3 (Convexidade da quadrática): Seja $f : \mathbb{R}^n \rightarrow \mathbb{R}$ uma função quadrática:

$$f(x) = x^T A x + 2b^T x + c \quad (81)$$

Onde A é simétrica. Então:

$$f \text{ (estritamente) convexa} \Leftrightarrow A \succeq 0 (A \succ 0) \quad (82)$$

Demonstração: A prova para o caso Pelo Teorema 2.1.1.1 e sabendo que $\nabla f(x) = 2(Ax + b)$, temos que f é convexa \Leftrightarrow

$$\forall x, y \in \mathbb{R}^n, y^T A y + 2b^T y + c \geq x^T A x + 2b^T x + c + 2(Ax + b)^T(y - x) \quad (83)$$

Rearranjando, obtemos:

$$\forall x, y \in \mathbb{R}^n (y - x)^T A (y - x) \geq 0 \Rightarrow A \succeq 0 \quad (84)$$

□

Teorema 2.1.1.4 (Monotonicidade do gradiente): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ continuamente diferenciável, então:

$$f \text{ convexa em } C \Leftrightarrow \forall x, y \in C, (\nabla f(x) - \nabla f(y))^T (x - y) \geq 0 \quad (85)$$

Demonstração: (\Rightarrow) Assuma que f é convexa sobre C . Por Teorema 2.1.1.1:

$$\begin{aligned} f(x) &\geq f(y) + \nabla f(y)^T (x - y) \\ f(y) &\geq f(x) + \nabla f(x)^T (y - x) \end{aligned} \quad (86)$$

Somando ambas as igualdades, obtemos (85)

(\Leftarrow) Suponha que (85) seja válida e sejam $x, y \in C$, vamos definir a função:

$$g(t) := f(x + t(y - x)), \quad t \in [0, 1] \quad (87)$$

Pelo Teorema Fundamental do Cálculo:

$$\begin{aligned} f(y) &= g(1) = g(0) + \int_0^1 g'(t) dt \\ &= f(x) + \int_0^1 (y - x)^T \nabla f(x + t(y - x)) dt \\ &= f(x) + (y - x)^T \nabla f(x) + \int_0^1 (y - x)^T (\nabla f(x + t(y - x)) - \nabla f(x)) dt \\ &= f(x) + (y - x)^T \nabla f(x) + \frac{1}{t} \int_0^1 t(y - x)^T (\nabla f(x + t(y - x)) - \nabla f(x)) dt \\ &\geq f(x) + (y - x)^T \nabla f(x) \end{aligned} \quad (88)$$

Onde utilizamos (85) na última desigualdade

□

2.1.2 Caracterizações de convexidade de segunda ordem

Teorema 2.1.2.1 (Caracterização de convexidade de segunda ordem): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ duas vezes continuamente diferenciável sobre um conjunto convexo C , então:

$$f \text{ convexa em } C \Leftrightarrow \forall x \in C, \nabla^2 f(x) \succeq 0 \quad (89)$$

Demonstração: (\Leftarrow) Suponha que $\nabla^2 f(x) \succeq 0$ para todo $x \in C$. Sejam $x, y \in C$. Pelo teorema de aproximação linear, existe $\xi \in [x, y] \subset C$ tal que

$$f(y) = f(x) + \nabla f(x)^T (y - x) + (y - x)^T \nabla^2 f(\xi) (y - x) \quad (90)$$

Como $\nabla^2 f(\xi) \succeq 0$, segue que

$$f(y) = f(x) + \nabla f(x)^T(y - x) \quad (91)$$

Como o argumento vale para todo $x, y \in C$, provamos que f é convexa em C pelo Teorema 2.1.1.1.

(\Rightarrow) Suponha que f é convexa em C . Sejam $x \in C$ e $d \in \mathbb{R}^n$ com $\|d\| = 1$. Sendo C aberto, existe $\varepsilon > 0$ tal que $x + \lambda d \in C$ para todo $0 < \lambda < \varepsilon$. Para tal λ , segue do Teorema 2.1.1.1

$$f(x + \lambda d) \geq f(x) + \lambda \nabla f(x)^T d \quad (92)$$

Além disso, pelo teorema de aproximação quadrática:

$$f(x + \lambda d) = f(x) + \lambda \nabla f(x)^T d + \frac{\lambda^2}{2} d^T \nabla^2 f(x) d + o(\lambda^2 \|d\|^2) \quad (93)$$

Combinando as expressões, obtemos, para todo $\lambda \in (0, \varepsilon)$

$$\frac{\lambda^2}{2} d^T \nabla^2 f(x) d + o(\lambda^2) \geq 0 \quad (94)$$

Isso é:

$$d^T \nabla^2 f(x) d + \frac{o(\lambda^2)}{\lambda^2} \geq 0 \quad (95)$$

Fazendo com que $\lambda \rightarrow 0$, temos:

$$d^T \nabla^2 f(x) d \geq 0 \quad (96)$$

Ou seja, $\nabla^2 f(x) \succeq 0, \forall x$ □

Teorema 2.1.2.2 (Caracterização de convexidade de segunda ordem): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ duas vezes continuamente diferenciável sobre um conjunto convexo C , então:

$$\forall x \in C, \nabla^2 f(x) \succ 0 \Rightarrow f \text{ estritamente convexa em } C \quad (97)$$

A volta na questão anterior não vale, por exemplo, $f(x) = x^4$ tem mínimo em 0, mas $f''(0) = 0$

2.1.3 Convexidade forte

Vimos o conceito de convexidade aplicando a condição de pontos numa reta estarem acima da curva da função. Mas e se uma forma mais curvada ainda tivesse em cima da função?

Definição 2.1.3.1 (Convexidade forte): Uma função $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ com C convexo é μ -fortemente convexa ($\mu > 0$) se:

$$\begin{aligned} &\forall x, y \in C \wedge \forall \lambda \in [0, 1], \\ &f(\lambda x + (1 - \lambda)y) + \frac{\mu}{2} \lambda(1 - \lambda) \|y - x\|^2 \leq \lambda f(x) + (1 - \lambda)f(y) \end{aligned} \quad (98)$$

Mas o que diabos isso significa? A gente agora, em vez de checar se os pontos na reta $t(x, f(x)) + (1 - t)(y, f(y))$ ($t \in [0, 1]$), imagine que tem uma cordinha entre esses dois pontos, a gravidade vai afetar ela e ela vai ficar curvada, e μ dita o quão curvada a cordinha está. Se os pontos nessa

cordinha estão acima da curva para todos os pontos na curva, então a função é μ -fortemente convexa

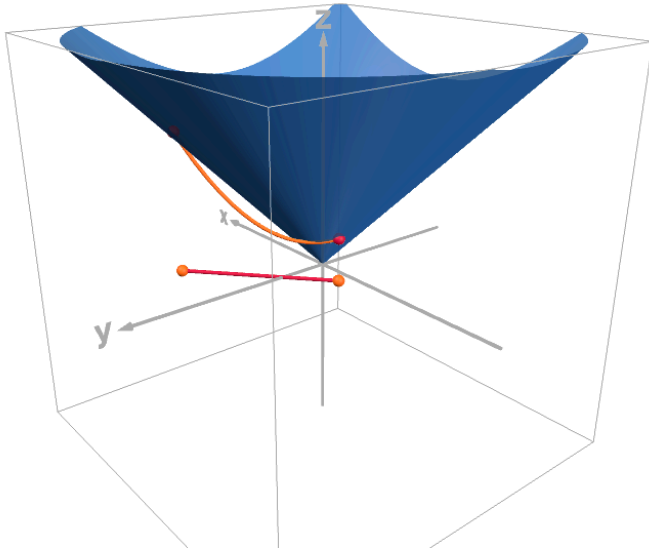


Figura 10: Função não-fortemente convexa, mas convexa $f(x) = \|x\|$

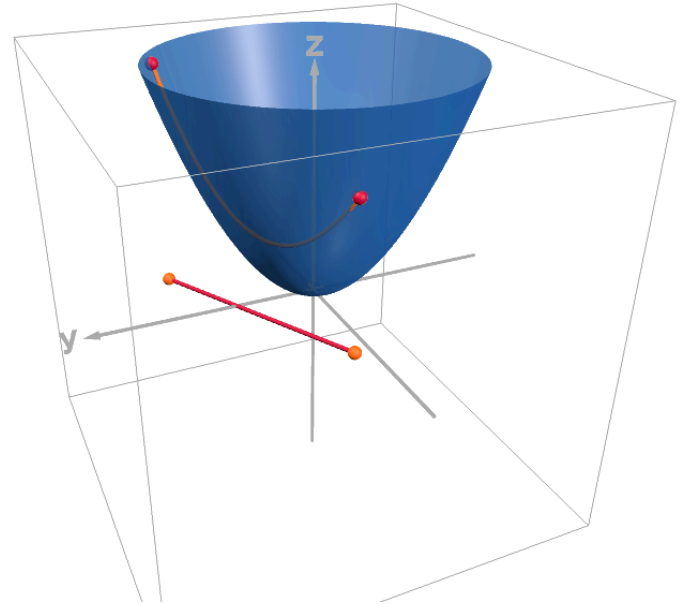


Figura 11: Função fortemente convexa $f(x) = \|x\|^2$

Teorema 2.1.3.1 (Desigualdade do gradiente: fortemente convexa): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ continuamente diferenciável e C convexo. Temos:

$$f \text{ } \mu\text{-fortemente convexa} \Leftrightarrow \forall x, y \in C, f(x) + \nabla f(x)^T(y - x) + \frac{\mu}{2} \|y - x\|^2 \leq f(y) \quad (99)$$

Teorema 2.1.3.2 (Caracterização de convexidade forte de segunda ordem): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ duas vezes continuamente diferenciável e C convexo. Então:

$$f \text{ } \mu\text{-fortemente convexa} \Leftrightarrow \nabla^2 f(x) - \mu I \succ 0 \quad (100)$$

2.2 Otimização sobre conjuntos convexos

Com toda essa bagagem, conseguimos finalmente aplicar a otimização de f em uma restrição convexa C

$$\min_{x \in C} f(x) \quad (101)$$

2.2.1 Condição de primeira ordem: Caso geral

Vamos primeiramente ver uma condição sobre funções generalizadas. Algo que faz sentido pensar quando estamos sendo restringidos, é pensar que não necessariamente meu máximo ou mínimo vai ter derivada igual a 0, veja o exemplo:

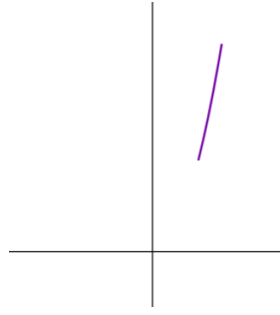


Figura 12: Exemplo de restrição: $f(x) = x^2$ com $x \in [2, 3]$

Teorema 2.2.1.1 (Condição de primeira ordem: Caso restrito): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ continuamente diferenciável em C convexo e fechado, então:

$$x^* \in C \text{ mínimo local} \Rightarrow \forall x \in C, \nabla f(x^*)^T (x - x^*) \geq 0 \quad (102)$$

Demonstração: Precisamos do seguinte lema: Seja $f : U \rightarrow \mathbb{R}$ função continuamente diferenciável sobre um aberto $U \subset \mathbb{R}^n$. Se para algum $x \in U$ e $d \neq 0$ tem-se

$$\nabla f(x)^T d < 0 \quad (103)$$

então existe $\varepsilon > 0$ tal que para todo $t \in (0, \varepsilon)$, $x + td \in U$ e

$$f(x + td) < f(x) \quad (104)$$

Continuemos a demonstração do teorema original. Assuma por contradição que exista $x \in C$ tal que $\nabla f(x^*)^T (x - x^*) < 0$. Temos então que, para $d := x - x^*$, $f'(x^*; d) = \nabla f(x^*)^T (x - x^*) < 0$. Segue do lema anterior, que existe $\varepsilon \in (0, 1)$ tal que

$$\forall t \in (0, \varepsilon), f(x^* + td) < f(x^*) \quad (105)$$

Sendo C convexo, segue que $x^* + td = (1 - t)x^* + tx \in C$. Concluimos então que x^* não é um ponto de mínimo local de f em C — uma contradição. \square

Mas o que esse teorema quer dizer??? Vamos por partes. Lembra do cosseno entre dois vetores v e u ?

$$\cos(\theta) = \frac{u^T v}{\|u\| \|v\|} \quad (106)$$

Ou seja, quando o sinal do ângulo entre eles depende única e exclusivamente de $u^T v$. Lembre que, se $\theta \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ então $\cos(\theta) \geq 0$ e se $\theta \in [\frac{\pi}{2}, \frac{3\pi}{2}]$ então $\cos(\theta) \leq 0$. Mas o que isso quer dizer? Espera mais um pouco. Lembra que vimos em cálculo 2 que o vetor gradiente indica a direção no domínio que eu devo seguir para que **a função aumente**? Show, agora a gente pode entender o que o teorema quer dizer para nós.

Vamos considerar o caso mais básico, quando x^* não ta na fronteira de C

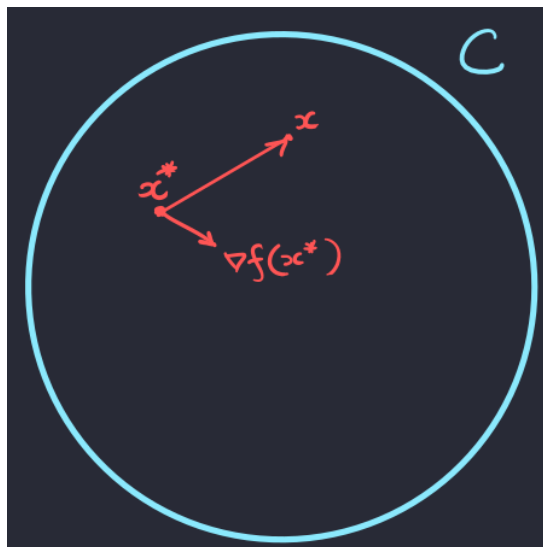


Figura 13: Ponto mínimo $x^* \in C$

Na imagem temos o vetor gradiente e o vetor $x - x^*$. Quando variamos o nosso ponto x , podemos claramente perceber que o vetor $x - x^*$ faz vários ângulos com o gradiente, só que se o gradiente for desça forma, ao andarmos na direção oposta ao gradiente, nossa função vai diminuir, ou seja, x^* não pode ser um ponto de mínimo! O que isso quer dizer? Que meu gradiente é 0!

Mas e se x^* estiver na minha fronteira?

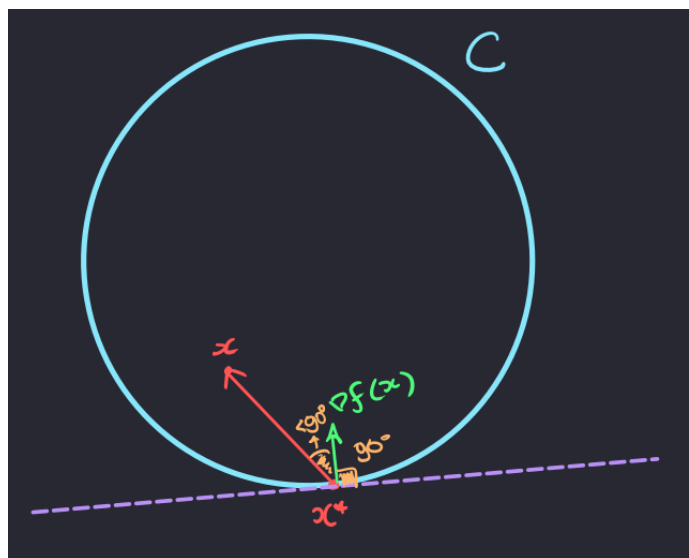


Figura 14: x^* mínimo na fronteira

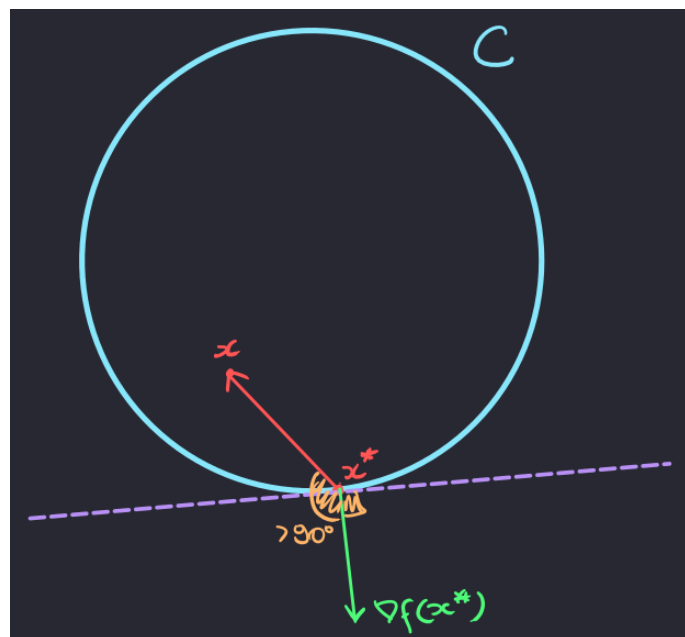


Figura 15: Função não-convexa $f(x) = x^4 - 3x^2$

Perceba que na primeira figura, se eu vejo o ângulo do gradiente com qualquer outro ponto no meu conjunto eu tenho menos que 90 graus, ou seja, o meu gradiente aponta para **dentro do conjunto**, de forma que a única maneira de diminuir mais a função é **saindo da restrição**. Na outra figura isso é melhor ilustrado. Veja que existem vetores no conjunto que fazem mais que 90 graus com o vetor gradiente, ou seja, o vetor gradiente ta para fora do conjunto C , de forma que eu consigo andar na direção $-\nabla^2 f(x^*)$ para que diminua ainda mais a função, ou seja, x^* não seria um mínimo

Esse teorema nos da motivação para uma definição

Definição 2.2.1.1 (Ponto estacionário): Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ com C convexo e fechado, chamamos $x^* \in C$ de ponto estacionário quando

$$\forall x \in C, \nabla f(x^*)(x - x^*) \geq 0 \quad (107)$$

2.2.2 Condições de primeira ordem: Caso convexo

Teorema 2.2.2.1: Seja $f : C \subset \mathbb{R}^n \rightarrow \mathbb{R}$ continuamente diferenciável e convexa com C convexo e fechado e $x^* \in C$, então:

$$x^* \text{ mínimo global} \Leftrightarrow x^* \text{ é ponto estacionário} \quad (108)$$

Demonstração: Precisamos provar apenas (\Leftarrow) do Teorema 2.2.1.1. Seja $x^* \in C$ um ponto estacionário de f em C . Obtemos que, para todo $x \in C$,

$$f(x) \geq f(x^*) + \nabla f(x^*)^T(x - x^*) \geq f(x^*) \quad (109)$$

onde a primeira desigualdade segue da desigualdade do gradiente (Teorema 2.1.1.1) e a segunda desigualdade segue de que x^* é ponto estacionário. Sendo que

$$\forall x \in C, f(x) \geq f(x^*) \quad (110)$$

segue que $x^* \in C$ é ponto de mínimo global de f em C . □

Otimização com restrições lineares

Aqui nós vamos introduzir um teorema muito importante no ramo da otimização, o **teorema das condições KKT**. Esse teorema generaliza as condições **necessárias** para um problema de minimização **genérico**, porém, vamos começar por baixo, em vez de já ir para o caso geral, vamos começar a passos pequenos

Primeiramente, queremos minimizar problemas do tipo:

$$\begin{aligned} \min_x f(x) \\ x \text{ sujeito a restrições do tipo } a_i^T x \leq b_i, \quad i = 1, \dots, m \end{aligned} \quad (111)$$

onde f é continuamente diferenciável em \mathbb{R}^n , $\{a_i\}_{i=1}^m \subset \mathbb{R}^n$, $\{b_i\}_{i=1}^m \subset \mathbb{R}$. Ou seja, o conjunto viável C é o políedro:

$$C = \cap_{i=1}^m \{x \in \mathbb{R}^n / a_i^T x \leq b_i\} \quad (112)$$

Há um exemplo nas anotações sobre convexidade do Phillip que mostram que C é convexo.

3.1 Condições KKT

Teorema 3.1.1 (Condições KKT para restrições lineares: condições necessárias de otimalidade): Considere o problema de minimização (111) onde f é uma função continuamente diferenciável em \mathbb{R}^n , $\{a_i\}_{i=1}^m \subset \mathbb{R}^n$ e $\{b_i\}_{i=1}^m \subset \mathbb{R}$. Então, **se** x^* é um ponto de **mínimo local** do problema, $\exists \lambda_1, \dots, \lambda_m \geq 0$ tais que

$$\begin{aligned} \nabla f(x^*) + \sum_{i=1}^m \lambda_i a_i &= 0, \\ \lambda_i (a_i^T x^* - b_i) &= 0, \quad i = 1, \dots, m \\ a_i^T x^* - b_i &\leq 0, \quad i = 1, \dots, m \end{aligned} \quad (113)$$

Como esse teorema necessita de vários outros resultados, não vou escrever a sua demonstração aqui. Se estiver curioso para saber a demonstração, confira o apêndice das anotações do Phillip

3.2 Condições KKT: Problema convexo

Teorema 3.2.1 (Condições KKT para restrições lineares: condições necessárias de otimalidade com função convexa): Considere o problema de minimização

$$\begin{aligned} \min_x f(x) \\ \text{sujeito à } a_i^T x \leq b_i, \quad i = 1, \dots, m \end{aligned} \quad (114)$$

onde f é uma função continuamente diferenciável **convexa** em \mathbb{R}^n , $\{a_i\}_{i=1}^m \subset \mathbb{R}^n \wedge \{b_i\}_{i=1}^m \subset \mathbb{R}$. Então, se x^* é um ponto de mínimo local do problema $\Leftrightarrow \exists \lambda_1, \dots, \lambda_m \geq 0$ tais que

$$\begin{aligned} \nabla f(x^*) + \sum_{i=1}^m \lambda_i a_i &= 0, \\ \lambda_i (a_i^T x^* - b_i) &= 0, \quad i = 1, \dots, m \\ a_i^T x^* - b_i &\leq 0, \quad i = 1, \dots, m \end{aligned} \quad (115)$$

Demonstração: (\implies) Segue do Teorema 3.1.1

(\impliedby) Definamos a função:

$$h(x) := f(x) + \sum_{i=1}^m \lambda_i (a_i^T x - b_i) \quad (116)$$

Temos que:

$$\nabla h(x^*) = \nabla f(x^*) + \sum_{i=1}^m \lambda_i a_i \quad (117)$$

Como h é convexa (Soma de funções convexas), segue que x^* é ponto mínimo de h em \mathbb{R}^n . Em particular, dado qualquer $x \in \mathbb{R}^n$ tal que:

$$a_i^T x \leq b_i, \quad i = 1, \dots, m \quad (118)$$

Tem-se que:

$$\begin{aligned} f(x^*) &= f(x^*) + \sum_{i=1}^m \lambda_i (a_i^T x^* - b_i) \\ &\leq f(x^*) + \sum_{i=1}^m \lambda_i (a_i^T x - b_i) \\ &\leq f(x) \end{aligned} \quad (119)$$

Na primeira equação utilizamos a segunda condição e na segunda desigualdade usamos o fato que $\lambda_i \geq 0$. Concluimos então que x^* é solução do sistema \square

3.3 Condições KKT com restrições lineares de igualdade

Show! Vimos as restrições afins de **desigualdade**, porém, em alguns casos, é possível que tenhamos restrições de igualdade:

$$\begin{aligned} &\min_x f(x) \\ &x \text{ sujeito a restrições do tipo :} \\ &a_i^T x \leq b_i, \quad i = 1, \dots, m \\ &c_j^T x = d_j, \quad j = 1, \dots, p \end{aligned} \quad (120)$$

onde f é continuamente diferenciável em \mathbb{R}^n , $\{a_i\}_{i=1}^m \subset \mathbb{R}^n$, $\{b_i\}_{i=1}^m \subset \mathbb{R}$, $\{c_j\}_{j=1}^p \subset \mathbb{R}^n$

Esse caso é o que costumamos aprender em cálculo dois como o **método de Lagrange**, porém vamos ver que esse método é **bem** mais geral do que viamos antes. Do problema que estabelecemos antes, segue um teorema bem parecido com Teorema 3.1.1

Teorema 3.3.1: Considere o problema (120), onde f é continuamente diferenciável em \mathbb{R}^n , $\{a_i\}_{i=1}^m \subset \mathbb{R}^n$, $\{b_i\}_{i=1}^m \subset \mathbb{R}$, $\{c_j\}_{j=1}^p \subset \mathbb{R}^n$. Então:

a) Se x^* é um ponto de mínimo local do problema, então existem $\lambda_1, \dots, \lambda_m \geq 0$ e $\mu_1, \dots, \mu_p \in \mathbb{R}$ tais que

$$\begin{aligned}\nabla f(x^*) + \sum_{i=1}^m \lambda_i a_i + \sum_{j=1}^p \mu_j c_j &= 0 \\ \lambda_i (a_i^T x^* - b_i) &= 0, \quad i = 1, \dots, m \\ a_i^T x^* - b_i &\leq 0, \quad i = 1, \dots, m \\ \mu_j (c_j^T x^* - d_j) &= 0, \quad j = 1, \dots, p\end{aligned}\tag{121}$$

b) Suponha adicionalmente que f é convexa, então x^* é um mínimo global do problema \Leftrightarrow existem $\lambda_1, \dots, \lambda_m \geq 0$ e $\mu_1, \dots, \mu_p \in \mathbb{R}$ tais que as condições (121) ainda valem

Demonstração: Primeiro demonstraremos o (a). Demonstrar essa parte é equivalente a resolver o problema:

$$\begin{aligned}\min_x & f(x) \\ x & \text{ sujeito a restrições do tipo } a_i^T x \leq b_i, \quad i = 1, \dots, m \\ & c_j^T x \leq d_j \wedge -c_j^T x \leq -d_j, \quad j = 1, \dots, p\end{aligned}\tag{122}$$

onde f é continuamente diferenciável em \mathbb{R}^n , $\{a_i\}_{i=1}^m \subset \mathbb{R}^n$, $\{b_i\}_{i=1}^m \subset \mathbb{R}$, $\{c_j\}_{j=1}^p \subset \mathbb{R}^n$

Sendo x^* uma solução do problema descrito anteriormente, pelo Teorema 3.1.1, temos:

$$\begin{aligned}\nabla f(x^*) + \sum_{i=1}^m \lambda_i a_i + \sum_{j=1}^p \mu_j^+ c_j - \sum_{j=1}^p \mu_j^- c_j &= 0 \\ \lambda_i (a_i^T x^* - b_i) &= 0 \\ \mu_j^+ (c_j^T x^* - d_j) &= 0 \\ \mu_j^- (-c_j^T x^* + d_j) &= 0\end{aligned}\tag{123}$$

Como x^* é viável, então as segundas e terceiras condições mencionadas na reformulação anterior são satisfeitas. Definindo então $\mu_j = \mu_j^+ - \mu_j^-$, então temos que $\sum_{j=1}^p \mu_j^+ c_j - \sum_{j=1}^p \mu_j^- c_j = \sum_{j=1}^p \mu_j c_j$. Então segue que as condições estabelecidas originalmente no teorema são satisfeitas

Para a demonstração de (b), Suponha que x^* viável e existem $\lambda_1, \dots, \lambda_m \geq 0$ e $\mu_1, \dots, \mu_p \in \mathbb{R}$ tais que as condições do teorema sejam satisfeitas. Defina

$$\mu_j^+ := (\mu_j)_+ = \max\{\mu_j, 0\}, \quad \mu_j^- := (\mu_j)_- = \max\{-\mu_j, 0\}\tag{124}$$

Como $\mu_j = \mu_j^+ - \mu_j^-$ e $c_j^T x^* - d_j = 0$ para $j \in [p]$, segue em particular que (123) é satisfeito. Sendo f convexa, segue do Teorema 3.2.1 que x^* é solução do problema reformulado e, em particular, do problema original do teorema \square

Otimização com restrições genéricas

Vimos minimizações para condições lineares e convexas, mas nem sempre isso acontece, muitas vezes temos conjuntos de restrições completamente genéricas. Nesse capítulo vamos ver como as condições KKT podem ser generalizadas para esses tipos de problemas (Vão perceber que elas ainda se aplicam ao que foi estabelecido anteriormente nas condições lineares). Vamos, então, redefinir o problema que tínhamos anteriormente:

$$\begin{aligned} \min_x f(x) \\ g_i(x) \leq 0, \quad \forall i = 1, \dots, m \\ h_j(x) = 0, \quad \forall j = 1, \dots, p \end{aligned} \quad (125)$$

4.1 Lagrangeano

O lagrangeando é uma função que será de grande importância, ela pode parecer meio confusa (Pois ela é), mas, a partir de agora, ela será nossa definição de “Derivar e igualar a 0”. Como assim? Sempre que queríamos minimizar/maximizar uma função, derivávamos e igualávamos a 0, só que vimos que, com restrições, isso não funciona mais, porém, essa função ainda se aplica (Com algumas ressalvas) no lagrangeano (Como veremos)

Definição 4.1.1 (Lagrangeano): O **Lagrangeano** associado à função f é a função $L : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}$ tal que:

$$L(x, \lambda, \mu) = f(x) + \lambda^T g(x) + \mu^T h(x) \quad (126)$$

Onde:

$$\lambda = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_m \end{pmatrix}, \quad \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix}, \quad g(x) = \begin{pmatrix} g_1(x) \\ \vdots \\ g_m(x) \end{pmatrix}, \quad h(x) = \begin{pmatrix} h_1(x) \\ \vdots \\ h_p(x) \end{pmatrix} \quad (127)$$

Teorema 4.1.1 (Gradiente Lagrangeano): Dado o Lagrangeano de uma função f , temos que o gradiente do lagrangeano **somente em relação a x** se da por:

$$\nabla_x L(x, \lambda, \mu) = \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla g_i(x) + \sum_{j=1}^p \mu_j \nabla h_j(x) \quad (128)$$

4.2 As generalizações do KKT

Agora queremos generalizar totalmente o KKT, então vamos aos poucos. Lembre que, até que falemos o contrário, estamos considerando o problema (125)

Antes de entrarmos diretamente no teorema KKT generalizado, vamos agora fazer uma definição que tem uma razão matemática, mas acaba por nos ajudar em alguns casos. Essa definição evita condições redundantes no nosso problema, já que elas podem acabar nos atrapalhando. Faremos um exemplo para mostrar essa ajuda

Definição 4.2.1 (Condições de qualificação de independência linear): Sejam $g_1, \dots, g_m : \mathbb{R}^n \rightarrow \mathbb{R}$ e $h_1, \dots, h_p : \mathbb{R}^n \rightarrow \mathbb{R}$ continuamente diferenciáveis e $x^* \in \mathbb{R}^n$, defina:

$$I(x^*) := \{i \in [m] : g_i(x^*) = 0\} \quad (129)$$

Dizemos que LICQ (Linear Independent Condition Qualification) é satisfeita em x^* para as funções g_1, \dots, g_m e h_1, \dots, h_p se

$$\{\nabla g_i(x^*) : i \in I(x^*)\} \cup \{\nabla h_j(x^*) : j \in [p]\} \text{ é linearmente independente} \quad (130)$$

Definição feita, vamos enunciar o novo teorema KKT

Teorema 4.2.1 (KKT): Se x^* é um ponto de mínimo local de $f(x)$ no problema (125) e a Definição 4.2.1 é satisfeita em x^* , isso implica que:

$$\begin{aligned} \exists \lambda_1, \dots, \lambda_m \geq 0, \exists \mu_1, \dots, \mu_p \in \mathbb{R} \\ \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) + \sum_{j=1}^p \mu_j \nabla h_j(x^*) = 0 \\ \lambda_i g_i(x^*) = 0 \quad i \in [m] \\ g_i(x^*) \leq 0 \quad i \in [m] \\ h_j(x^*) = 0 \quad j \in [p] \end{aligned} \quad (131)$$

Exemplo (Utilidade da LICQ): aaaaaaaaaa preencher aqui aaaaaaaaaaaaaa

Agora que vimos esses teoremas e condições, vamos fazer uma definição para facilitar em algumas terminologias:

Definição 4.2.2 (Ponto KKT): Considere o problema (125), onde $f, g_1, \dots, g_m, h_1, \dots, h_p$ são continuamente diferenciáveis no \mathbb{R}^n . Um ponto x^* viável, ou seja, que satisfaz as condições do conjunto viável, é chamado de **ponto KKT** quando $\exists \lambda_1, \dots, \lambda_m \geq 0$ e $\exists \mu_1, \dots, \mu_p \in \mathbb{R}$ tais que:

$$\begin{aligned} \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) + \sum_{j=1}^p \mu_j \nabla h_j(x^*) = 0 \\ \lambda_i g_i(x^*) = 0 \quad i \in [m] \end{aligned} \quad (132)$$

Isso facilita um pouco a terminologia pois podemos resumir o Teorema 4.2.1 em dizer que um ponto de LICQ não pode ser um ponto de mínimo se ele não for KKT

4.3 Caso Convexo

Claro, não poderíamos de falar do caso convexo aqui, sempre tem algo de especial nele, vamos então enunciar novamente o nosso problema mudando ele um pouco

$$\begin{aligned}
& \min_x f(x) \\
& g_i(x) \leq 0, \quad \forall i = 1, \dots, m \\
& h_j(x) = 0, \quad \forall j = 1, \dots, p
\end{aligned} \tag{133}$$

onde f , g_i e h_j são convexas $\forall i$ e $\forall j$

Vale ressaltar também que h_j são **afins**

Teorema 4.3.1 (KKT Convexo): Se x^* é um ponto de mínimo local de f dado o problema (133) e a Definição 4.2.1 é satisfeita em x^* , então x^* é uma solução do problema se, e somente se:

$$\begin{aligned}
& \exists \lambda_1, \dots, \lambda_m \geq 0, \quad \exists \mu_1, \dots, \mu_p \in \mathbb{R} \\
& \nabla f(x^*) + \sum_{i=1}^m \lambda_i \nabla g_i(x^*) + \sum_{j=1}^p \mu_j \nabla h_j(x^*) = 0 \\
& \lambda_i g_i(x^*) = 0 \quad i \in [m] \\
& g_i(x^*) \leq 0 \quad i \in [m] \\
& h_j(x^*) = 0 \quad j \in [p]
\end{aligned} \tag{134}$$

Show! Inclusive, por conta que o nosso problema é convexo, podemos trocar a necessidade do LICQ (Definição 4.2.1) por uma condição um pouco mais fácil

Definição 4.3.1 (Condição de Slater): Dizemos que a condição de Slater é satisfeita para as funções g_1, \dots, g_m (convexas) se

$$\exists \hat{x} \in \mathbb{R}^n / g_i(\hat{x}) < 0, \quad \forall i \in [m] \tag{135}$$

Ou seja, essa condição é satisfeita quando x^* é um ponto viável (Dentro do conjunto viável). Agora podemos refazer o teorema utilizando dessa condição

Teorema 4.3.2 (KKT e Slater): Se x^* é mínimo local de $f(x)$ (Nas restrições $g_i(x) \leq 0$ e $h_j(x) = 0$ sendo funções continuamente diferenciáveis e convexas) e x^* satisfaz Definição 4.3.1, então x^* é ponto KKT (A volta não vale)

Porém, como falei anteriormente, não faz sentido falarmos de funções convexas de igualdade ($h_j(x) = 0$) que **não são afins**, isso nos permite reescrever o problema de uma forma interessante:

$$\begin{aligned}
& \min_x f(x) \\
& g_i(x) \leq 0 \quad i \in [m] \\
& h_j(x) = 0 \quad j \in [p] \\
& s_k(x) = 0 \quad k \in [q]
\end{aligned} \tag{136}$$

Onde f , g_i são convexas e h_j , s_k são afins

Então podemos adaptar a condição de Slater:

Teorema 4.3.3 (Condição de Slater): Dizemos que a condição de Slater é satisfeita para as funções g_1, \dots, g_m (convexas) e h_1, \dots, h_p e s_1, \dots, s_q (afins) quando:

$$\begin{aligned} \exists \hat{x} \in \mathbb{R}^n \text{ tal que} \\ g_i(\hat{x}) < 0, \quad \forall i \in [m] \\ h_j(\hat{x}) \leq 0, \quad \forall j \in [p] \\ s_k(\hat{x}) = 0, \quad \forall k \in [q] \end{aligned} \tag{137}$$

De forma que o Teorema 4.3.2 continua valendo. Vimos 3 tipos de condições diferentes! Que tal a gente refazer o nosso teorema de uma forma geral?

Teorema 4.3.4 (Final KKT): Dado o problema (125), e seja $x^* \in C$ (Conjunto viável), temos 3 caracterizações:

1. As restrições em x^* são LICQ
2. Problema convexo (133) + Condição de Slater (Definição 4.3.1)

Se x^* ou o problema satisfaz qualquer uma dessas condições, então eu posso dividir meu problema em algumas condições:

- (Necessidade com qualificação de restrições)
 - Se x^* é um mínimo local e as restrições ativas em x^* satisfazem LICQ $\Rightarrow x^*$ é um ponto KKT
- (Convexidade + Slater)
 - x^* é KKT (Ser mínimo \Rightarrow KKT)
 - Reciprocamente, se x^* é viável e é KKT, então x^* é ótimo global

Algoritmos de Otimização

Agora que vimos bastante da teoria por trás da otimização, precisamos viabilizar isso para problemas enormes! Não faz sentido, por exemplo, um ser humano comum resolver um sistema KKT para uma função no \mathbb{R}^{1000} por exemplo, né? Então vamos utilizar dos **computadores**. Temos algumas famílias de métodos para otimização:

- Métodos de ordem-zero
- Métodos de primeira-ordem
 - Método do gradiente (Cauchy)
 - Método do subgradiente
 - ...
- Método de segunda-ordem
 - Método de Newton
 - Métodos Quasi-Newton
 - Método de pontos-interiores
 - ...
- Método de confiança

Esse são apenas alguns exemplos, mas existem **inúmeros** algoritmos nesse ramo

5.1 Método Gradiente

Vamos definir o algoritmo do gradiente

```

1  $x_1$  inicial
2 for  $t = 1, \dots, n$  do:
3    $| x_{t+1} = x_t - \alpha_t \nabla f(x_t)$ 
```

Onde $\alpha_t > 0$ é o “passo” ou learning rate. Vamos agora fazer algumas definições para mostrar o porquê do método do gradiente funcionar

Definição 5.1.1 (Suavidade): Dizemos que $f : \mathbb{R}^n \rightarrow \mathbb{R}$ é uma função L -suave se $\exists L > 0$ tal que:

$$\forall x, y, \quad \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad (138)$$

Ou, equivalentemente

$$\|\nabla f(x) - \nabla f(y)\| \leq O(\|x - y\|) \quad (139)$$

Definição 5.1.2 (Direção de Descida): Dizemos que $d \in \mathbb{R}^n$ é de “descida” a partir de um ponto $x \in \mathbb{R}^n$ se:

$$(\nabla f(x))^T d < 0 \quad (140)$$

Teorema 5.1.1 (Suavidade): $\forall x, y \in \mathbb{R}^n$ vale que:

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \text{ é } L\text{-suave} \Leftrightarrow f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{L}{2} \|y - x\|^2 \quad (141)$$

Se d é a direção de descida em x , então

$$x^+ := x + \alpha d \quad (142)$$

tem que, por suavidade de f :

$$\begin{aligned} f(x^+) - f(x) &\leq \nabla f(x)^T (x^+ - x) + \frac{L}{2} \|x^+ - x\|^2 \\ &= \alpha \nabla f(x)^T d + \frac{L\alpha^2}{2} \|d\|^2 \\ &= \alpha \left(\nabla f(x)^T d + \frac{L\alpha}{2} \|d\|^2 \right) < 0 \end{aligned} \quad (143)$$

E temos que

$$\begin{aligned} \lim_{\alpha \rightarrow 0^+} \left(\nabla f(x)^T d + \frac{L\alpha}{2} \|d\|^2 \right) &= \nabla f(x)^T d < 0 \\ \Rightarrow \exists \hat{\alpha} > 0 \text{ tal que } \nabla f(x)^T d + \frac{L\hat{\alpha}}{2} \|d\|^2 &< 0 \end{aligned} \quad (144)$$

Teorema 5.1.2: $\forall t \in \mathbb{N}$ e supondo que $0 < \alpha_t \leq \frac{2}{L}$, então (Considerando que a direção de descida é $-\nabla f(x)$):

$$\alpha_t \left(1 - \frac{L\alpha_t}{2} \right) \|\nabla f(x_t)\|^2 \leq f(x_t) - f(x_{t+1}) \quad (145)$$

Demonstração: f é suave:

$$\begin{aligned} f(x_{t+1}) &\leq f(x_t) + \nabla f(x_t)^T (x_{t+1} - x_t) + \frac{L}{2} \|x_{t+1} - x_t\|^2 \\ &= f(x_t) - \alpha_t \|\nabla f(x_t)\|^2 + \frac{L\alpha_t^2}{2} \|\nabla f(x_t)\|^2 \end{aligned} \quad (146)$$

□