

Introduction to Machine Learning

Lecture 7 - Introduction to Optimization

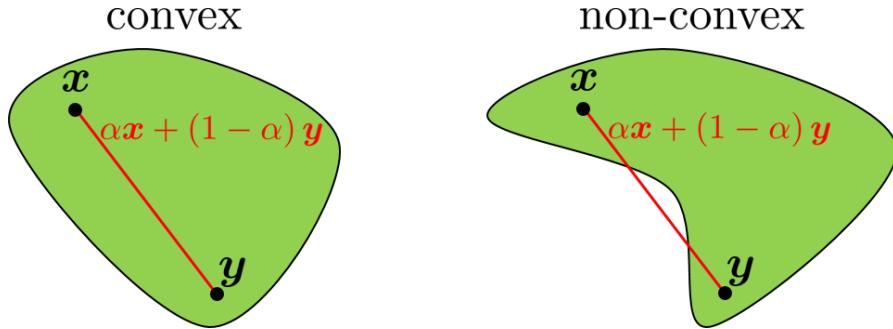
1 Convexity

1.1 Convex set

A set $\mathcal{C} \subseteq \mathbb{R}^d$ is said to be **convex** if, for all $\mathbf{x}, \mathbf{y} \in \mathcal{C}$ and for all $\alpha \in [0, 1]$:

$$\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in \mathcal{C}$$

For example:



Exercise 1 Let \mathcal{C} be a convex set and let $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3 \in \mathcal{C}$.

Prove that any **convex combination** is also in \mathcal{C} . namely:

$$\sum_{i=1}^3 \alpha_i \mathbf{x}_i \in \mathcal{C}$$

where $\alpha_i \geq 0$ and $\sum_{i=1}^3 \alpha_i = 1$.

Solution:

$$\begin{aligned} \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \alpha_3 \mathbf{x}_3 &= \underbrace{(\alpha_1 + \alpha_2)}_{1-\alpha_3} \underbrace{\left(\frac{\alpha_1}{\alpha_1 + \alpha_2} \mathbf{x}_1 + \frac{\alpha_2}{\alpha_1 + \alpha_2} \mathbf{x}_2 \right)}_{\triangleq \mathbf{y} \in \mathcal{C}} + \alpha_3 \mathbf{x}_3 \\ &= (1 - \alpha_3) \mathbf{y} + \alpha_3 \mathbf{x}_3 \in \mathcal{C} \end{aligned}$$

■

Exercise 2 Let \mathcal{C} and \mathcal{D} be convex sets.

Prove that $\mathcal{C} \cap \mathcal{D}$ is also convex set.

Solution:

Consider $\mathbf{x}, \mathbf{y} \in \mathcal{C} \cap \mathcal{D}$.

Since both \mathcal{C} and \mathcal{D} are convex, for any $\alpha \in [0, 1]$ we have:

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in \mathcal{C}, \quad \mathbf{x}, \mathbf{y} \in \mathcal{C}$$

And also:

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in \mathcal{D}, \quad \mathbf{x}, \mathbf{y} \in \mathcal{D}$$

Thus:

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in \mathcal{C} \cap \mathcal{D}, \quad \mathbf{x}, \mathbf{y} \in \mathcal{C} \cap \mathcal{D}$$

■

1.2 Convex function

Let \mathcal{C} be a convex set and let $f : \mathcal{C} \rightarrow \mathbb{R}$ be a function.

- f is called **convex** if, for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}$ and for all $\alpha \in [0, 1]$:

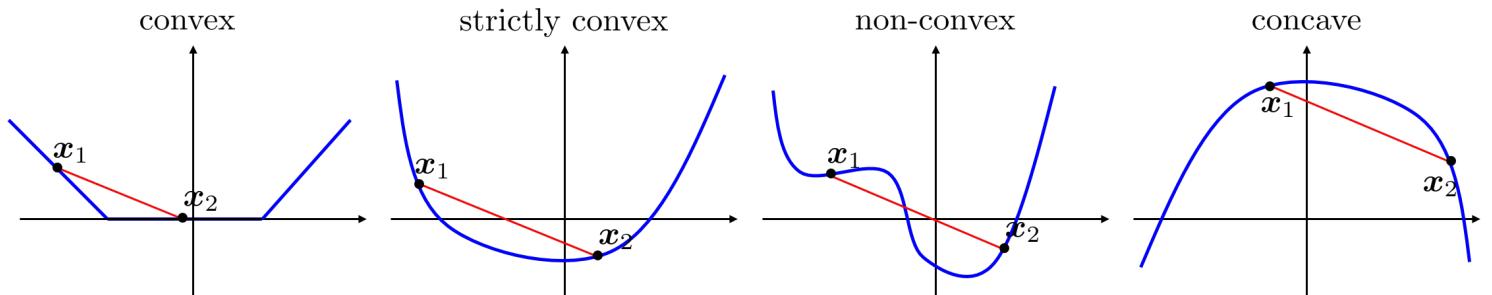
$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$$

- f is called **strictly convex** if, for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}$ such that $\mathbf{x}_1 \neq \mathbf{x}_2$ and for all $\alpha \in (0, 1)$:

$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) < \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$$

- f is called (strictly) concave if $-f$ is (strictly) convex.

For example:



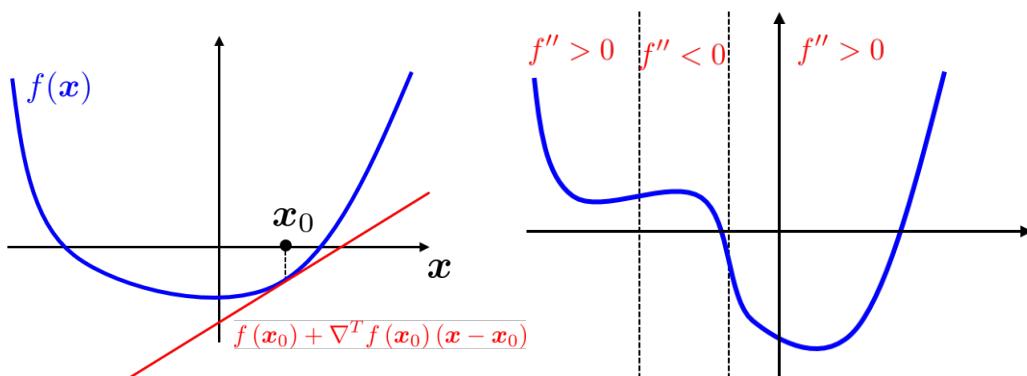
Theorem 1. Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is twice differentiable.

Then the following are equivalent:

1. f is convex.
2. $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla^T f(\mathbf{x})(\mathbf{y} - \mathbf{x})$ for all \mathbf{x} and \mathbf{y} .
3. $\nabla^2 f(\mathbf{x}) \succeq 0$ for all \mathbf{x} .

Condition 2 states that the Taylor expansion at any point is a global under estimator of f .

Condition 3 states that the function f has non-negative curvature everywhere.



f is strictly convex if, for all \mathbf{x} : $\nabla^2 f(\mathbf{x}) \succ 0$ (it is enough but it is not a necessary condition).

We prove the theorem for dimension one: $n = 1$.

1 \Rightarrow 2 Since f is convex, by definition for $\alpha \in (0, 1)$:

$$\begin{aligned} f(\alpha\mathbf{y} + (1 - \alpha)\mathbf{x}) &\leq \alpha f(\mathbf{y}) + (1 - \alpha) f(\mathbf{x}) \\ f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) &\leq f(\mathbf{x}) + \alpha(f(\mathbf{y}) - f(\mathbf{x})) \\ \frac{f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\alpha} &\leq f(\mathbf{y}) - f(\mathbf{x}) \end{aligned}$$

For $\alpha \rightarrow 0^+$ we have:

$$f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x})$$

$$\Rightarrow \boxed{f(\mathbf{y}) \geq f(\mathbf{x}) + f'(\mathbf{x})(\mathbf{y} - \mathbf{x})}$$

2 \Rightarrow 1 We have that:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) \quad \forall \mathbf{x}, \mathbf{y}$$

Consider:

$$\mathbf{z} \triangleq \alpha\mathbf{y} + (1 - \alpha)\mathbf{x}$$

$$\Rightarrow \begin{cases} f(\mathbf{y}) \geq f(\mathbf{z}) + f'(\mathbf{z})(\mathbf{y} - \mathbf{z}) & (1) \\ f(\mathbf{x}) \geq f(\mathbf{z}) + f'(\mathbf{z})(\mathbf{x} - \mathbf{z}) & (2) \end{cases}$$

$$\begin{aligned} \alpha(1) + (1 - \alpha)(2) &\implies \alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{x}) \geq f(\mathbf{z}) + f'(\mathbf{z})(\alpha(\mathbf{y} - \mathbf{z}) + (1 - \alpha)(\mathbf{x} - \mathbf{z})) \\ \alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{x}) &\geq f(\mathbf{z}) + f'(\mathbf{z}) \left(\underbrace{\alpha\mathbf{y} + (1 - \alpha)\mathbf{x} - \mathbf{z}}_{= \mathbf{z}} \right) \\ \alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{x}) &\geq f(\mathbf{z}) \end{aligned}$$

$$\Rightarrow \boxed{\alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{x}) \geq f(\alpha\mathbf{y} + (1 - \alpha)\mathbf{x})}$$

2 \Rightarrow 3 Consider $y > x$. From 2 we have:

$$\begin{aligned} \begin{cases} f(y) \geq f(x) + f'(x)(y - x) \\ f(x) \geq f(y) + f'(y)(x - y) \end{cases} &\Rightarrow \begin{cases} f(y) - f(x) \geq f'(x)(y - x) \\ f(y) - f(x) \leq f'(y)(y - x) \end{cases} \\ &\Rightarrow f'(y)(y - x) \geq f'(x)(y - x) \\ f'(y)(y - x) - f'(x)(y - x) &\geq 0 \\ \frac{f'(y) - f'(x)}{y - x} &\geq 0 \end{aligned}$$

For $y \rightarrow x^+$, we have:

$$\Rightarrow \boxed{f''(x) \geq 0}$$

3 \Rightarrow 2 Using Taylor series (with the mean-value form of the remainder):

$$f(y) = f(x) + f'(x)(y - x) + \frac{1}{2}f''(z)(y - x), \quad z \in [x, y]$$

$$\Rightarrow \boxed{f(y) \geq f(x) + f'(x)(y - x)}$$

1.2.1 Results

Consider the convex function $f : \mathcal{C} \rightarrow \mathbb{R}$:

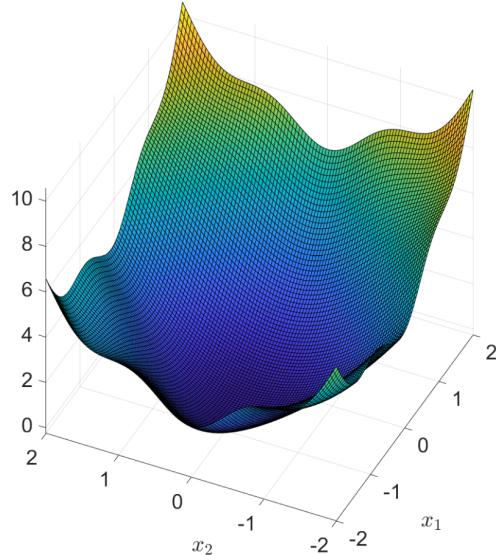
1. Any local minimum point is also a global minimum point.
2. The level set $\mathcal{D}_c = \{\mathbf{x} : f(\mathbf{x}) \leq c\}$ is a convex set.
3. If f is strictly convex, then the (global) minimum point (if exist) is unique.
4. If f is differentiable and $\mathbf{x}^* \in \mathcal{C}$ is not on the boundary of \mathcal{C} ($\mathbf{x}^* \notin \partial\mathcal{C}$) then:

$$\mathbf{x}^* \text{ is a minimum point} \Leftrightarrow \nabla f(\mathbf{x}^*) = \mathbf{0}$$

2 Unconstrained Optimization (gradient descent)

Consider the following function:

$$f(x_1, x_2) = x_1^2 + x_2^2 + \sin^2(x_1 x_2) + x_1$$



Find the minimum point of f :

$$\{x_1^*, x_2^*\} = \arg \min_{x_1, x_2} f(x_1, x_2) = ?$$

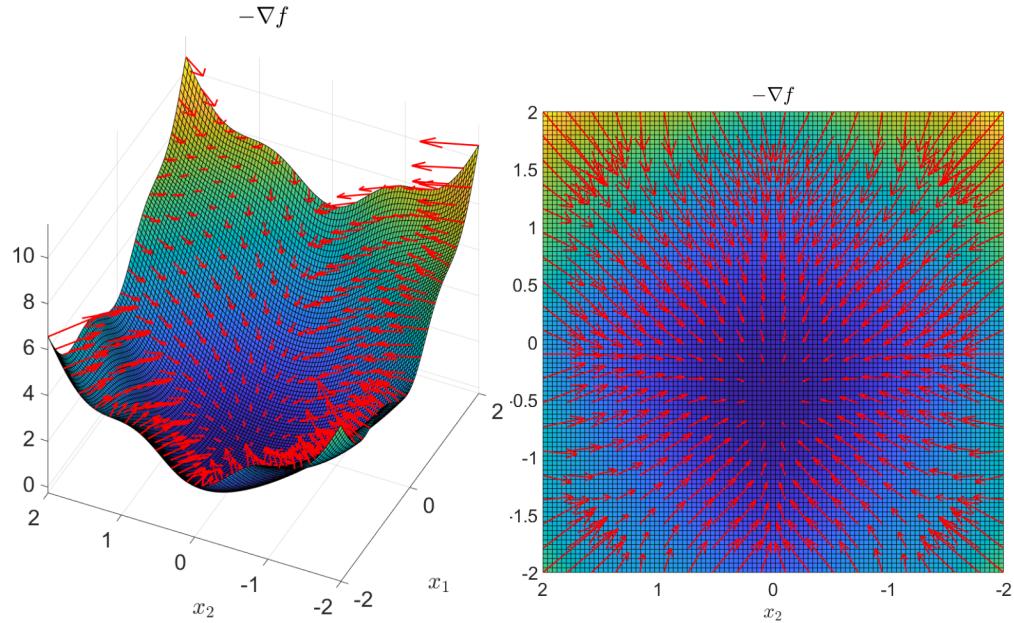
We can try to compare the gradient to zero:

$$\begin{aligned} \nabla f &= \mathbf{0} \\ \left[\begin{array}{l} \frac{\partial}{\partial x_1} f(x_1, x_2) \\ \frac{\partial}{\partial x_2} f(x_1, x_2) \end{array} \right] &= \mathbf{0} \\ \left[\begin{array}{l} 2x_1 + x_2 \sin(2x_1 x_2) + 1 \\ 2x_2 + x_1 \sin(2x_1 x_2) \end{array} \right] &= \mathbf{0} \\ \Rightarrow \left\{ \begin{array}{l} 2x_1 + x_2 \sin(2x_1 x_2) + 1 = 0 \\ 2x_2 + x_1 \sin(2x_1 x_2) = 0 \end{array} \right. \end{aligned}$$

In general, this equations has no analytical solution.

2.1 Gradient descent

Consider the (minus) gradient of f at each point:



At each point, the gradient ∇f points in the direction of (locally) maximum increase in f . Thus, to find a minimum point of f one can “go” in the opposite direction: $-\nabla f$.

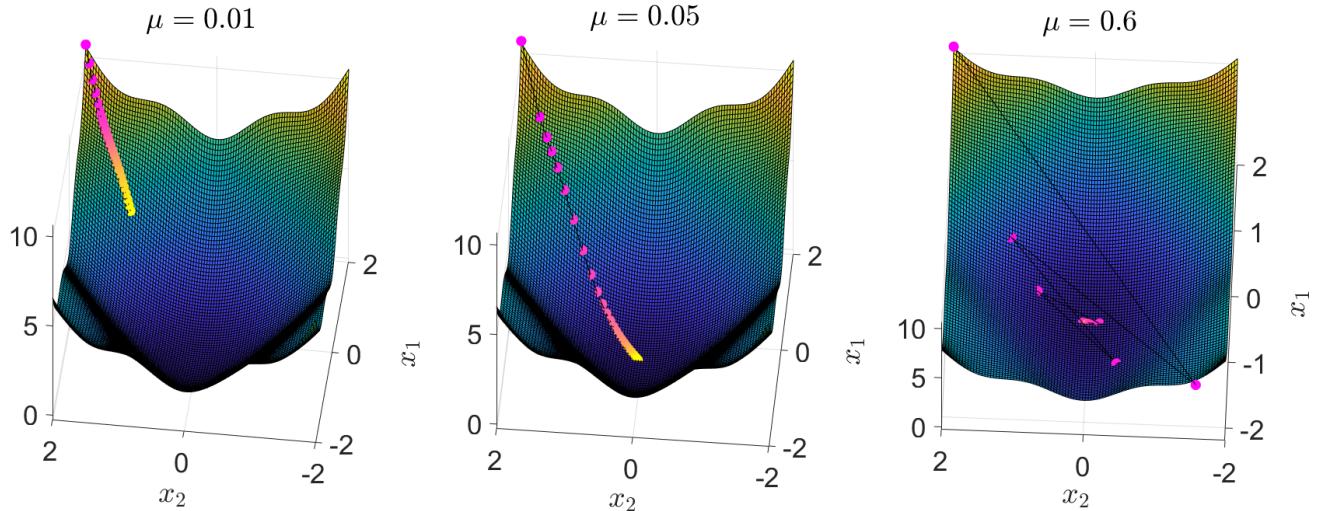
Algorithm 1 Gradient Descent Algorithm

1. Set some initial point \mathbf{x}_0
2. **for** $k = 0, 1, 2, \dots$ Update:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu \nabla f(\mathbf{x}_k)$$

The value of μ control the size of the step (learning rate).

We set $\mathbf{x}_0 = [2 \ 2]$ and apply the algorithm for 30 iterations:



2.2 Matrix inversion example

Consider the following 256×256 image:



We denote the column stack of the image with $\mathbf{x} \in \mathbb{R}^{65,536}$

Let \mathbf{H} be a blurring matrix and let \mathbf{y} be the blurred version (in column stack) of the image:

$$\mathbf{y} = \mathbf{H}\mathbf{x}, \quad \mathbf{H} \in \mathbb{R}^{65,536 \times 65,536}$$



To restore the original image \mathbf{x} from \mathbf{y} we need to compute:

$$\hat{\mathbf{x}} = \mathbf{H}^{-1}\mathbf{y}, \quad \text{Assuming } \mathbf{H} \text{ is invertible}$$

But, most computers cannot invert a $65,536 \times 65,536$ matrix.

Instead, we can defined the following optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} f(\mathbf{x}) = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2$$

The optimal solution is indeed $\hat{\mathbf{x}} = \mathbf{H}^{-1}\mathbf{y}$.

The gradient is given by:

$$\nabla f(\mathbf{x}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{x})$$

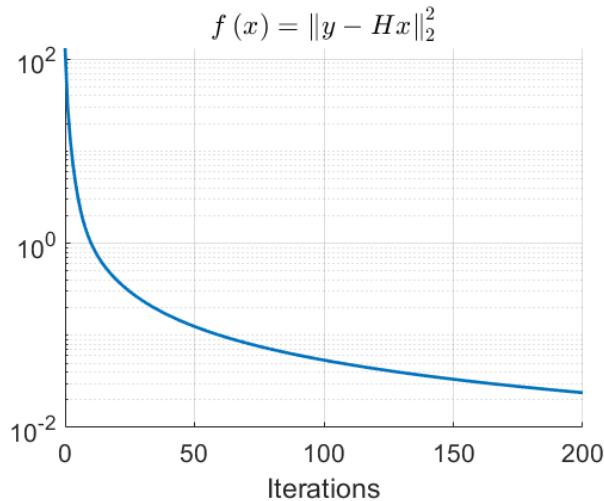
We set the initial solution:

$$\mathbf{x}_0 = \mathbf{y}$$

and $\mu = 0.9$.

We apply 200 iterations of the gradient descent algorithm.

We obtain values of $f(\mathbf{x}_k)$ (in logarithmic scale):



Some images during the iterations:



Final result:

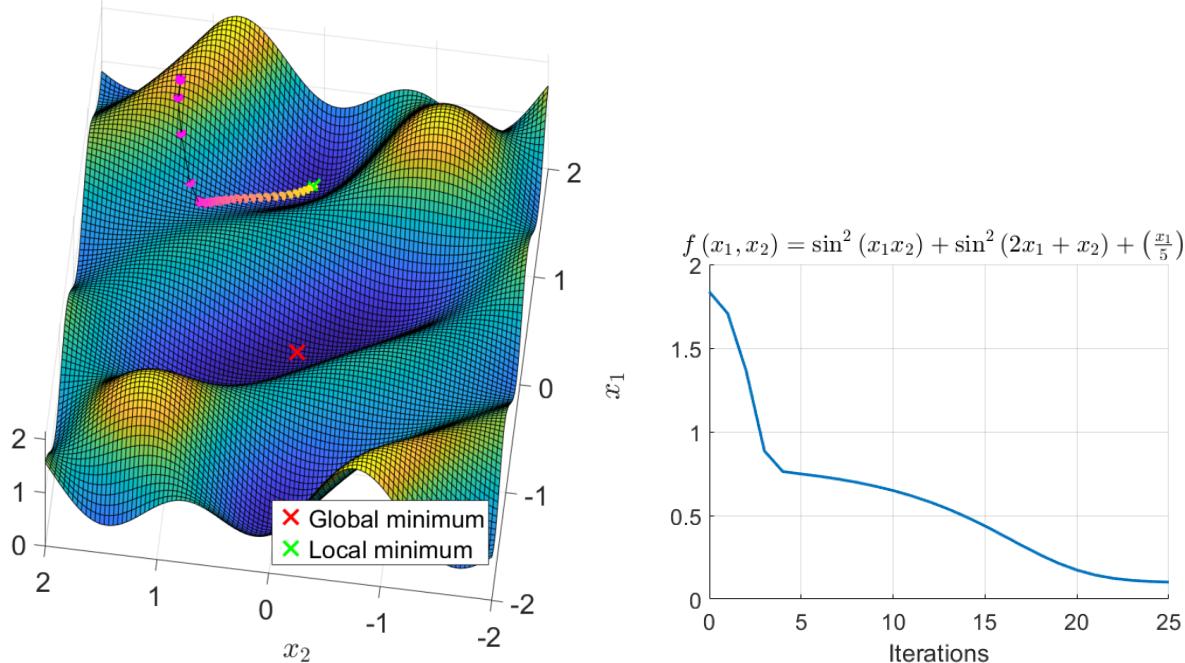


2.3 Non-convex function

- If the function f is convex, then, under some conditions on the step size μ , the gradient descent algorithm converges to the global minimum.
- If the function f is non-convex, then, the gradient might only converge to a local minimum.

Example:

$$f(x_1, x_2) = \sin^2(x_1 x_2) + \sin^2(2x_1 + x_2) + \left(\frac{x_1}{5}\right)^2$$



The gradient descent algorithm converge to a local minimum (note that $f(0, 0) = 0$).

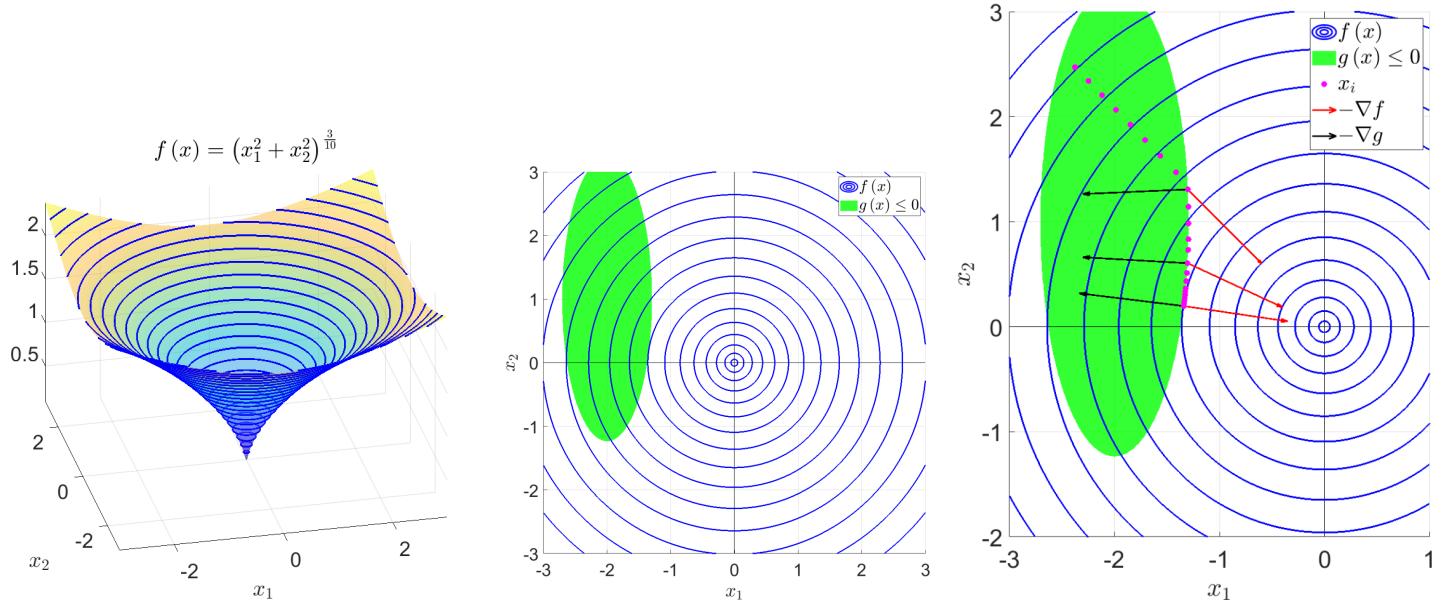
3 Constrained Optimization

Consider the following constrained optimization problem:

$$\begin{cases} \arg \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{subject to} \\ g(\mathbf{x}) \leq 0 \end{cases} = \begin{cases} \arg \min_{\mathbf{x}} (x_1^2 + x_2^2)^{\frac{3}{10}} \\ \text{subject to} \\ 2(x_1 + 2)^2 + \frac{1}{5}(x_2 - 1)^2 - 1 \leq 0 \end{cases}$$

Consider a starting point \mathbf{x}_0 inside the feasible area: $g(\mathbf{x}_0) \leq 0$.

Applying gradient descend to f while satisfying the condition $g \leq 0$ yields the following iterations (right figure):



Note that the algorithm converged to a point such that:

$$\nabla f(\mathbf{x}^*) = -\lambda \nabla g(\mathbf{x}^*), \quad \text{for some } \lambda > 0$$

$$\Rightarrow \boxed{\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0}$$

Thus, we write the Lagrangian (λ is called Lagrangian multiplier):

$$\boxed{\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})} \implies \boxed{\nabla_{\mathbf{x}} \mathcal{L} = \nabla f(\mathbf{x}) + \lambda \nabla g(\mathbf{x})}$$

The optimal point \mathbf{x}^* satisfying the KKT conditions:

1. Optimality under the constraint: $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}^*, \lambda^*) = 0$
2. Lagrange multiplier is non-negative: $\lambda^* \geq 0$
3. Solution is in the feasible area: $g(\mathbf{x}^*) \leq 0$
4. Complimentary slackness: $\lambda^* g(\mathbf{x}^*) = 0$

Condition 4 states that at least $g(\mathbf{x}^*) = 0$ (active constraint) or $\lambda^* = 0$ (non-active constraint).

3.1 Examples

3.1.1 Active constraint

Solve the following constraint problem:

$$\begin{cases} \arg \min_{\mathbf{x}} \|\mathbf{x}\|_2^2 \\ \text{subject to} \\ (x_1 + 2)^2 + (x_2)^2 - 1 \leq 0 \end{cases} \quad \mathbf{x} \in \mathbb{R}^2$$

Note that $\mathbf{x} = \mathbf{0}$ is the minimizer of $\|\mathbf{x}\|^2$ but it is not satisfy the constraint.

Solution:

The Lagrangian is given by ($\|\mathbf{x}\|^2 = x_1^2 + x_2^2$):

$$\mathcal{L} = x_1^2 + x_2^2 + \lambda ((x_1 + 2)^2 + (x_2)^2 - 1)$$

$$\Rightarrow \nabla_{\mathbf{x}} \mathcal{L} = \mathbf{0}$$

$$\begin{bmatrix} 2x_1 + \lambda(2x_1 + 4) \\ 2x_2 + 2\lambda x_2 \end{bmatrix} = \mathbf{0}$$

$$\Rightarrow \mathbf{x}^* \triangleq \begin{bmatrix} x_1^* \\ x_2^* \end{bmatrix} = \begin{bmatrix} -\frac{2\lambda}{1+\lambda} \\ 0 \end{bmatrix}$$

Complimentary slackness:

$$\lambda g(\mathbf{x}^*) = 0$$

If $\lambda = 0$ then $\mathbf{x}^* = \mathbf{0}$ but $g(\mathbf{0}) > 0$, namely, $\mathbf{x} = \mathbf{0}$ is not in the feasible area. Thus, $\lambda \neq 0$ (g is an active constraint):

$$\Rightarrow g(\mathbf{x}^*) = 0$$

$$(x_1^* + 2)^2 + (x_2^*)^2 - 1 = 0$$

$$(x_1^* + 2)^2 = 1$$

$$-\frac{2\lambda}{1+\lambda} + 2 = \pm 1$$

$$\begin{aligned} -\frac{2\lambda}{1+\lambda} + 2 &= -1 \\ 2\lambda &= 3(1+\lambda) \\ \lambda &= -3 \end{aligned}$$

$$\begin{aligned} -\frac{2\lambda}{1+\lambda} + 2 &= 1 \\ 2\lambda &= 1+\lambda \\ \lambda &= 1 \end{aligned}$$

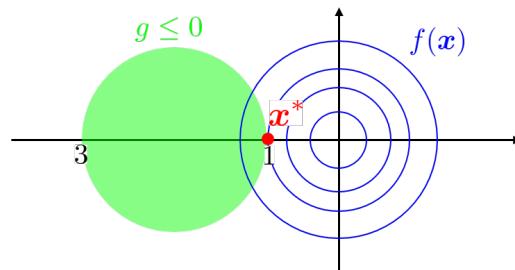
The Lagrange multiplier need to be non-negative: $\lambda \geq 0$

This is result is valid.

$$\lambda = -3$$

$$\lambda = 1$$

$$\Rightarrow \mathbf{x}^* = \begin{bmatrix} x_1^* \\ x_2^* \end{bmatrix} = \begin{bmatrix} -\frac{2\lambda}{1+\lambda} \\ 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$



3.1.2 Non-active constraint

Solve the following constraint problem:

$$\begin{cases} \arg \min_{\mathbf{x}} \|\mathbf{x}\|_2^2 \\ \text{subject to} \\ (x_1 + 2)^2 + (x_2)^2 - 10 \leq 0 \end{cases} \quad \mathbf{x} \in \mathbb{R}^2$$

Solution:

As before, the Lagrangian is given by ($\|\mathbf{x}\|^2 = x_1^2 + x_2^2$):

$$\mathcal{L} = x_1^2 + x_2^2 + \lambda ((x_1 + 2)^2 + (x_2)^2 - 1)$$

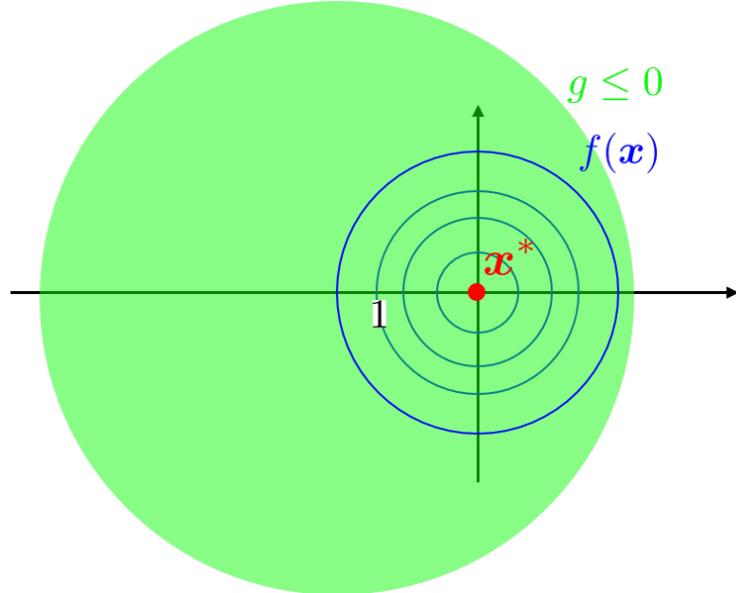
$$\begin{aligned} & \Rightarrow \nabla_{\mathbf{x}} \mathcal{L} = \mathbf{0} \\ & \begin{bmatrix} 2x_1 + \lambda(2x_1 + 4) \\ 2x_2 + 2\lambda x_2 \end{bmatrix} = \mathbf{0} \\ & \Rightarrow \mathbf{x}^* \triangleq \begin{bmatrix} x_1^* \\ x_2^* \end{bmatrix} = \begin{bmatrix} -\frac{2\lambda}{1+\lambda} \\ 0 \end{bmatrix} \end{aligned}$$

Complimentary slackness:

$$\lambda g(\mathbf{x}^*) = 0$$

If $\lambda = 0$ then $\mathbf{x}^* = \mathbf{0}$ and also $g(\mathbf{0}) < 0$, namely, $\mathbf{x} = \mathbf{0}$ is within the feasible area. Thus, $\lambda = 0$ and the constraint is not active:

$$\Rightarrow [\mathbf{x}^* = \mathbf{0}]$$



3.2 General Case

A general constrained optimization problem is given by:

$$\begin{cases} \arg \min_{\boldsymbol{x}} f(\boldsymbol{x}) \\ \text{subject to} \\ g_i(\boldsymbol{x}) \leq 0 & i = 1, 2, \dots, N \\ h_j(\boldsymbol{x}) = 0 & j = 1, 2, \dots, M \end{cases}$$

g_i are the inequality constraints and h_j are the equality constraints.

The Lagrangian is given by:

$$\mathcal{L} = f(\boldsymbol{x}) + \sum_i \lambda_i g_i(\boldsymbol{x}) + \sum_j \mu_j h_j(\boldsymbol{x})$$

If the functions $\{f, \{g_i\}, \{h_j\}\}$ are convex,

the point \boldsymbol{x} is an optimal solution if the generalized KKT (Karus-Khun-Tucker) conditions are satisfied.