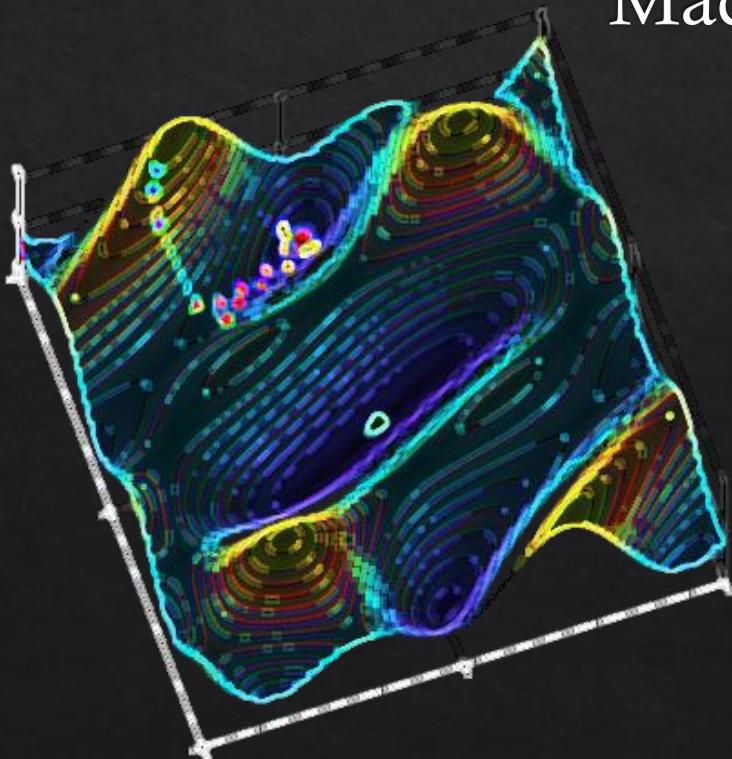


Introduction to Convex Optimization

Machine Learning Methods – Lecture 2



June 2021



Or Yair

The $\arg \min(\cdot)$ Operator – 1D

- Consider the quadratic function $f : \mathbb{R} \rightarrow \mathbb{R}$:

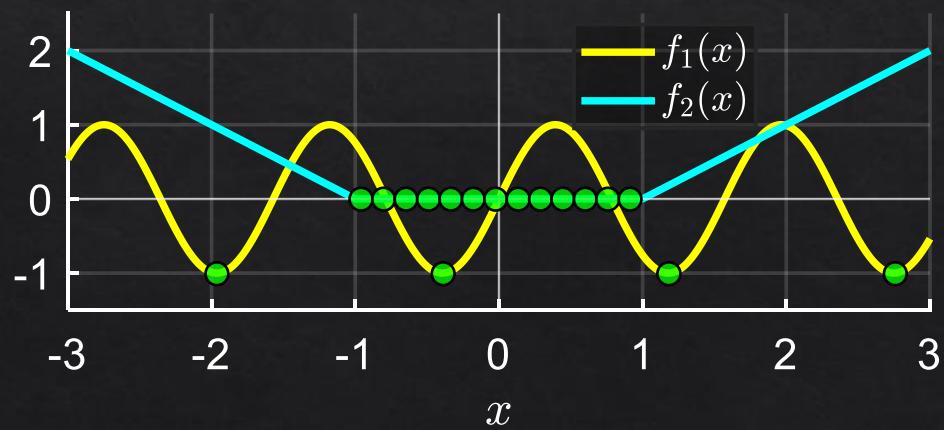
$$f(x) = ax^2 + bx + c, \quad a > 0$$

Find: $x^* = \arg \min_{x \in \mathbb{R}} f(x)$

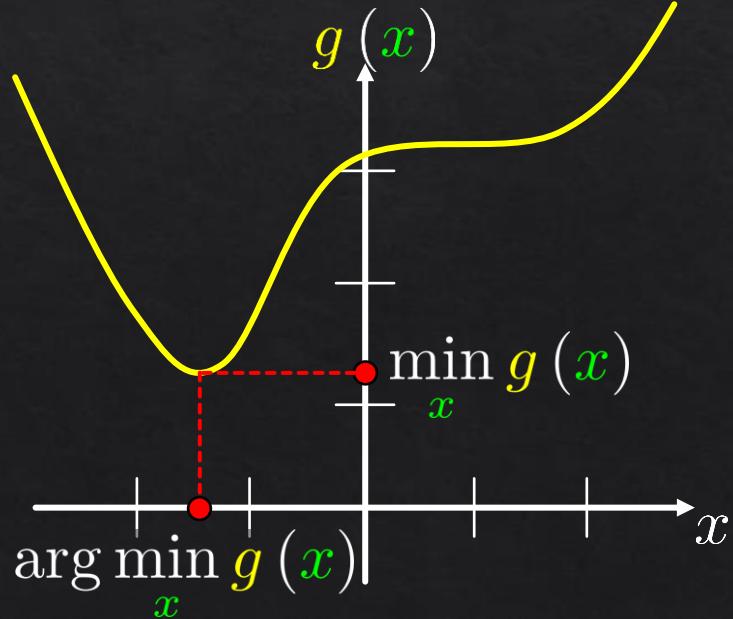
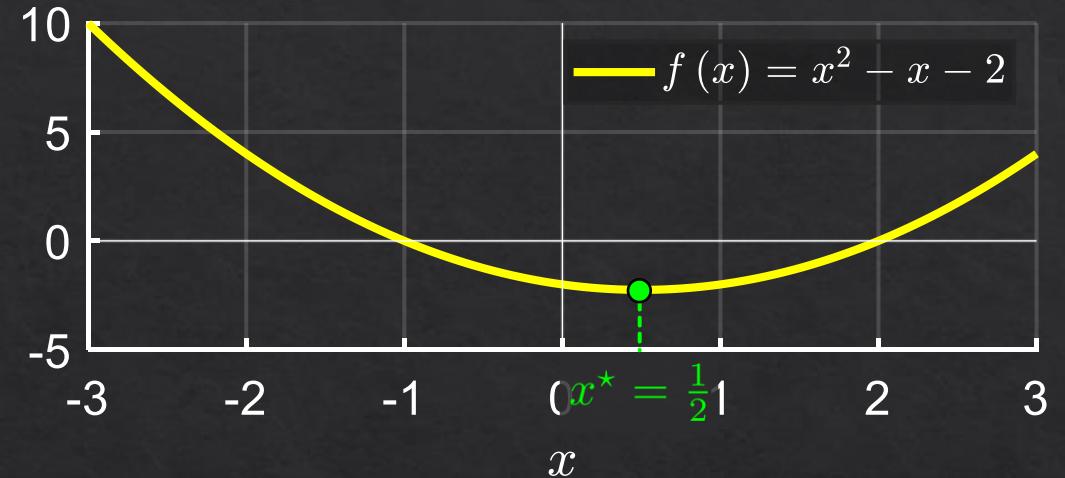
Solution

We can compare the derivative to zero:

$$\begin{aligned} f'(x^*) &= 0 \\ 2ax^* + b &= 0 \implies x^* = -\frac{b}{2a} \end{aligned}$$



- $\arg \min(\cdot)$ is not always well-defined (not unique).



The $\arg \min(\cdot)$ Operator – 2D

- Consider the quadratic function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$:

$$f(x_1, x_2) = a_1 x_1^2 + a_2 x_2^2 + b_1 x_1 + b_2 x_2 + c, \quad a_1, a_2 > 0$$

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

Solution:

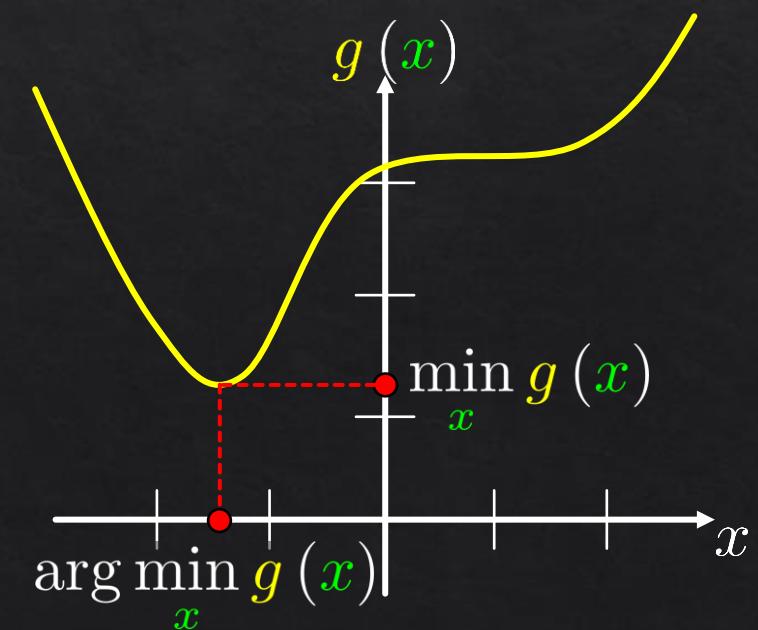
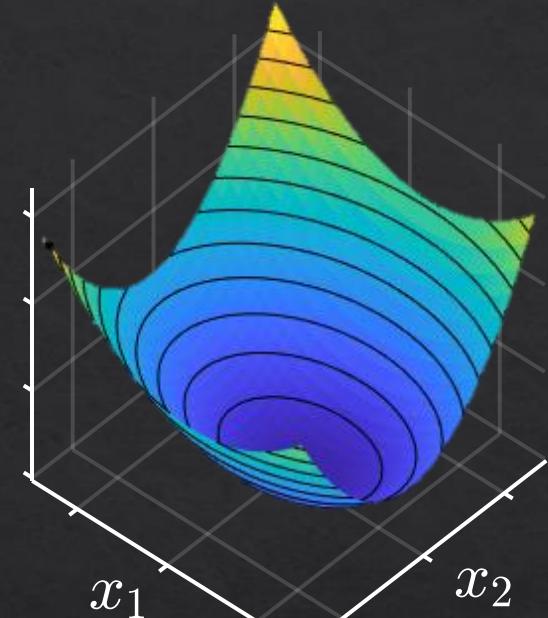
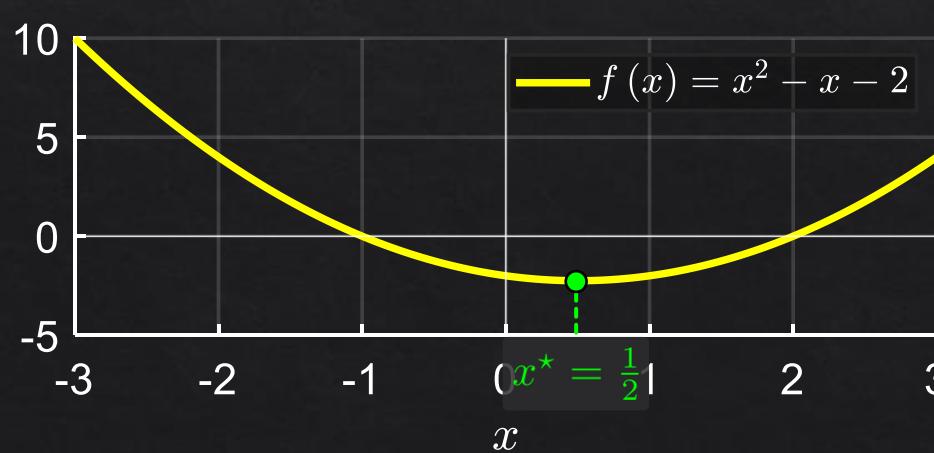
$$\begin{bmatrix} \frac{\partial f(x_1^*, x_2^*)}{\partial x_1} \\ \frac{\partial f(x_1^*, x_2^*)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 2a_1 x_1^* + b_1 \\ 2a_2 x_2^* + b_2 \end{bmatrix} = 0$$

$$\Rightarrow \begin{bmatrix} x_1^* \\ x_2^* \end{bmatrix} = \begin{bmatrix} -\frac{b_1}{2a_1} \\ -\frac{b_2}{2a_2} \end{bmatrix}$$

$$\nabla f(x_1, x_2) := \begin{bmatrix} \frac{\partial f(x_1, x_2)}{\partial x_1} \\ \frac{\partial f(x_1, x_2)}{\partial x_2} \end{bmatrix}$$

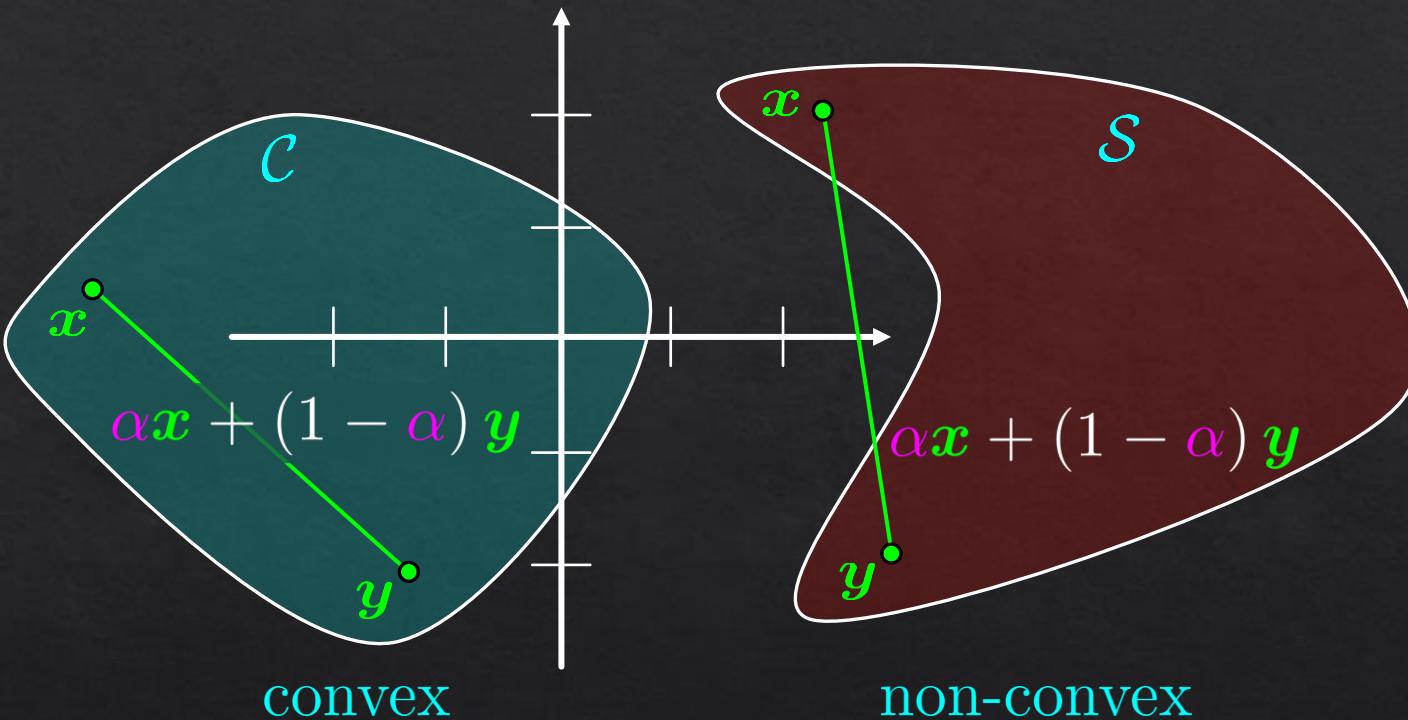
the gradient of f



Convex Set

- A set $\mathcal{C} \subseteq \mathbb{R}^d$ is convex if, for all $\mathbf{x}, \mathbf{y} \in \mathcal{C}$:

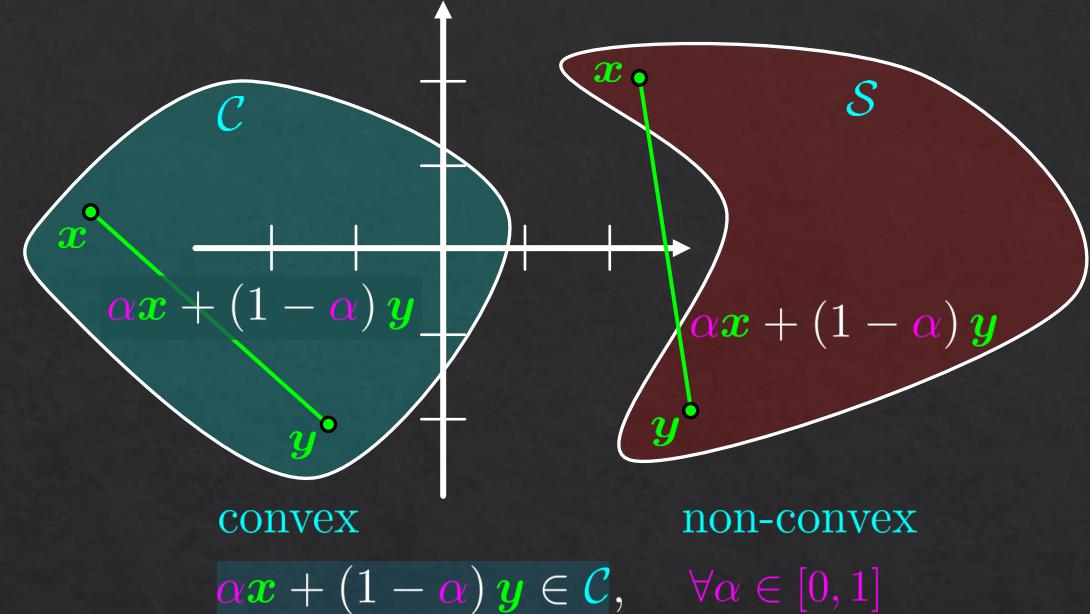
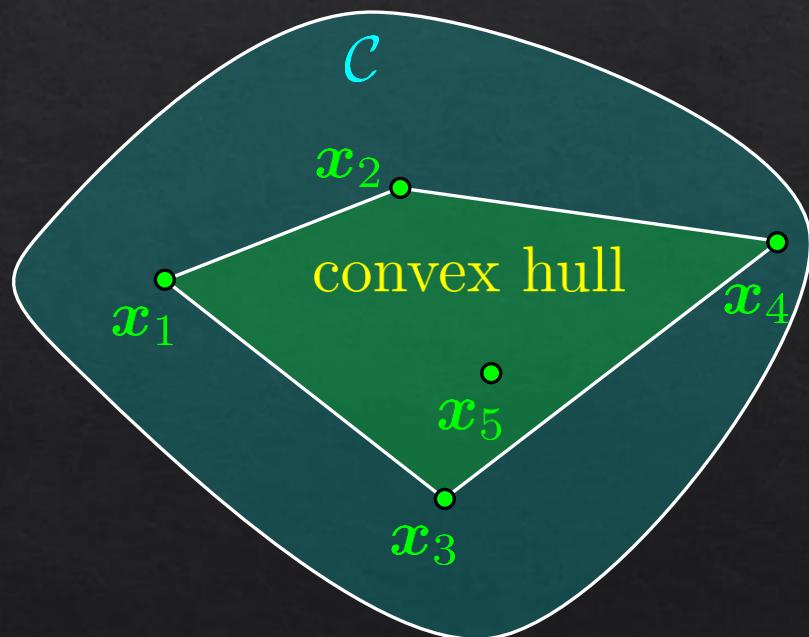
$$\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in \mathcal{C}, \quad \forall \alpha \in [0, 1]$$



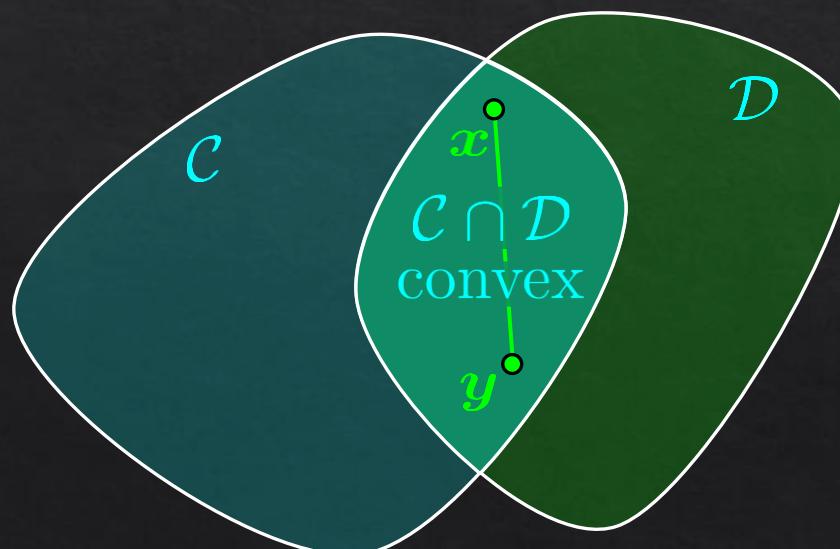
Convex Set – Properties

- Consider $\mathcal{X} = \{\mathbf{x}_i \in \mathcal{C}\}_{i=1}^N$ where \mathcal{C} is convex.
- Any convex combination of \mathcal{X} is also in \mathcal{C} .

$$\sum_{i=1}^N \alpha_i \mathbf{x}_i \in \mathcal{C}, \quad \sum_{i=1}^N \alpha_i = 1, \quad \alpha_i \geq 0$$



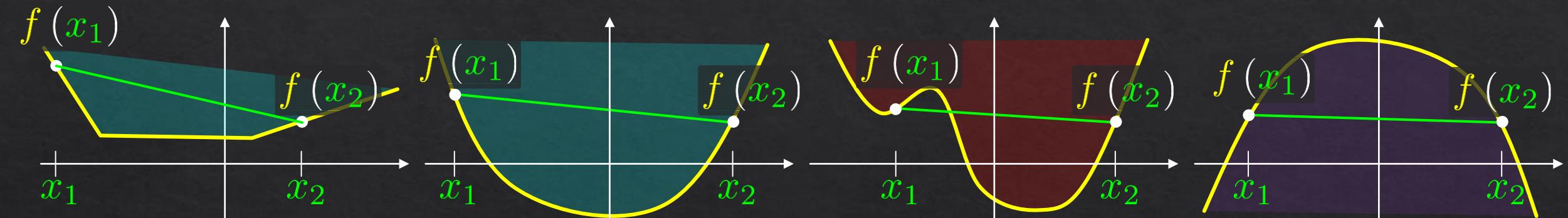
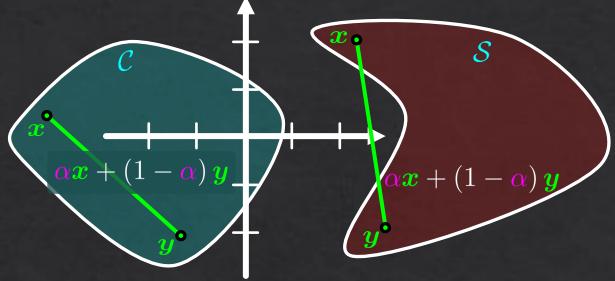
- If \mathcal{C} and \mathcal{D} are convex, so is $\mathcal{C} \cap \mathcal{D}$.



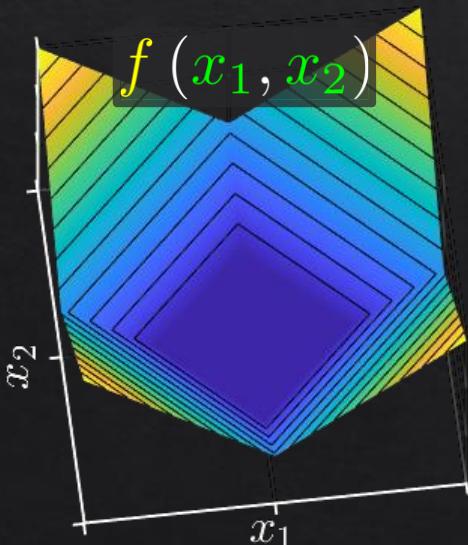
Convex Functions

- Consider $f : \mathcal{C} \rightarrow \mathbb{R}$ where \mathcal{C} is convex.
- f is called convex if, for all $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}$:

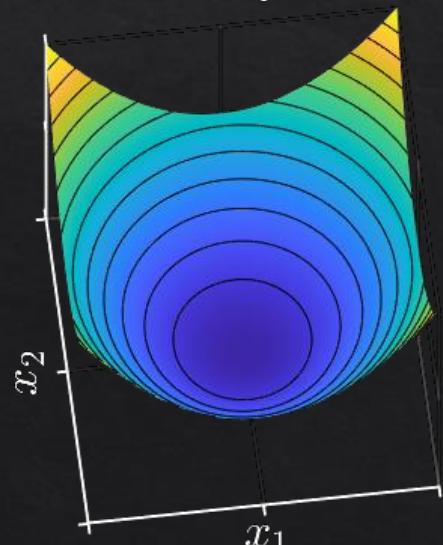
$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2), \quad \forall \alpha \in [0, 1]$$



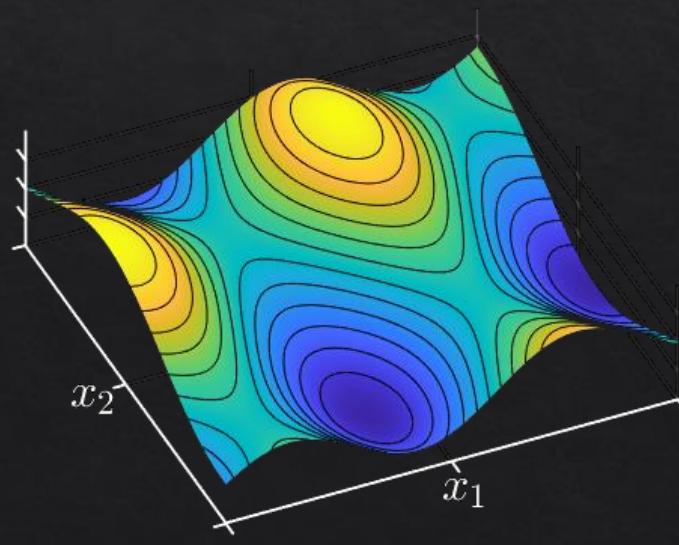
convex



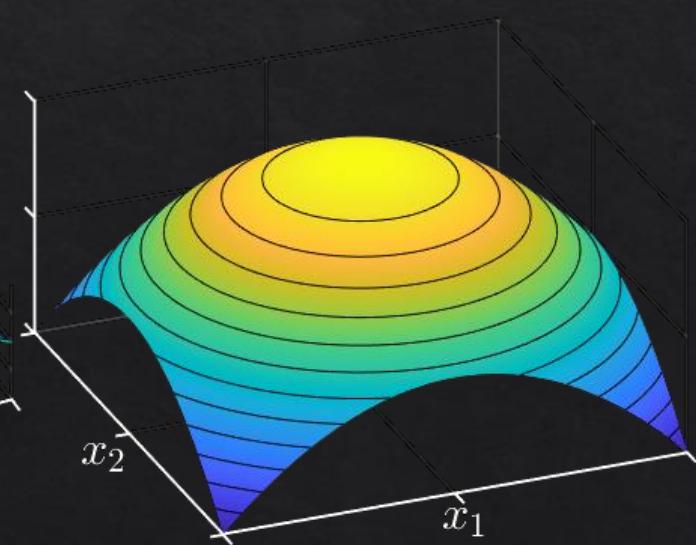
strictly convex



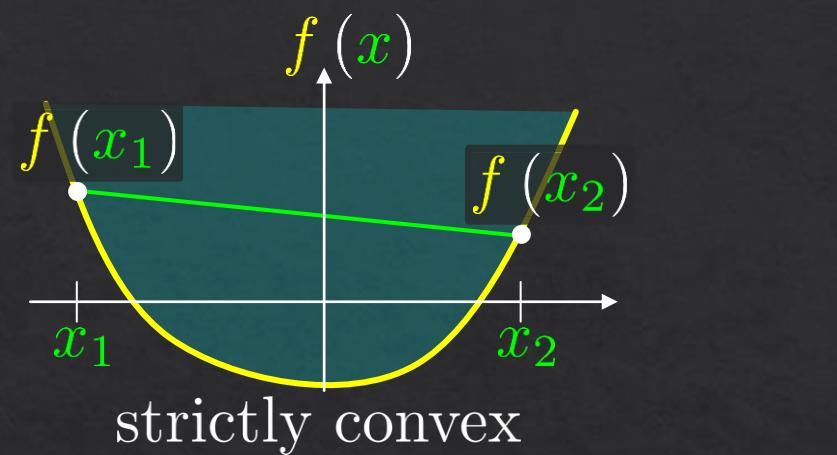
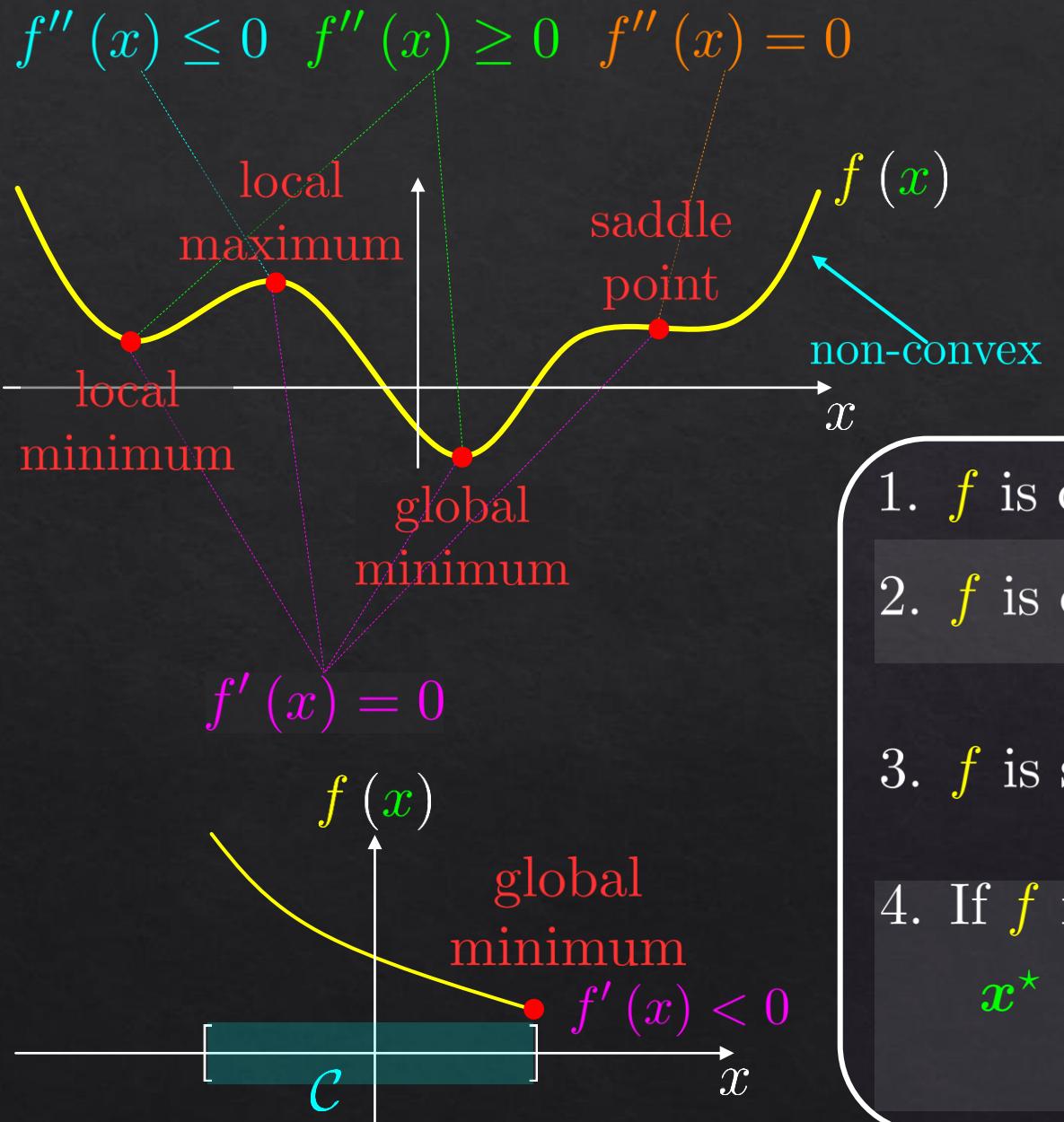
non-convex



concave*



Convex Functions – Properties



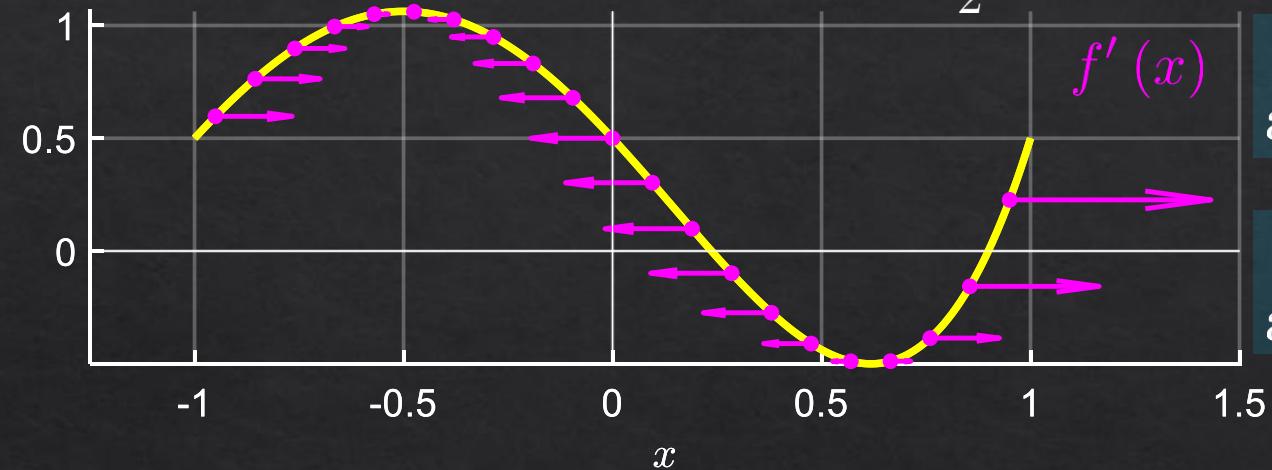
$$f(\alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2) \leq \alpha f(\mathbf{x}_1) + (1 - \alpha) f(\mathbf{x}_2)$$

1. f is convex $\iff \forall \mathbf{x} : f''(\mathbf{x}) \geq 0 \quad (\nabla^2 f(\mathbf{x}) \succeq 0)$
2. f is convex \implies local minimum = global minimum
3. f is strictly convex \implies the (global) minimum (if exist) is unique.
4. If f is convex, differentiable and $\mathbf{x}^* \notin \partial \mathcal{C}$, then:
 \mathbf{x}^* is a (global) minimum $\iff \nabla f(\mathbf{x}^*) = \mathbf{0}$
(f has no maxima)

Derivative (1D Gradient)

- Consider the following function $f : \mathbb{R} \rightarrow \mathbb{R}$:

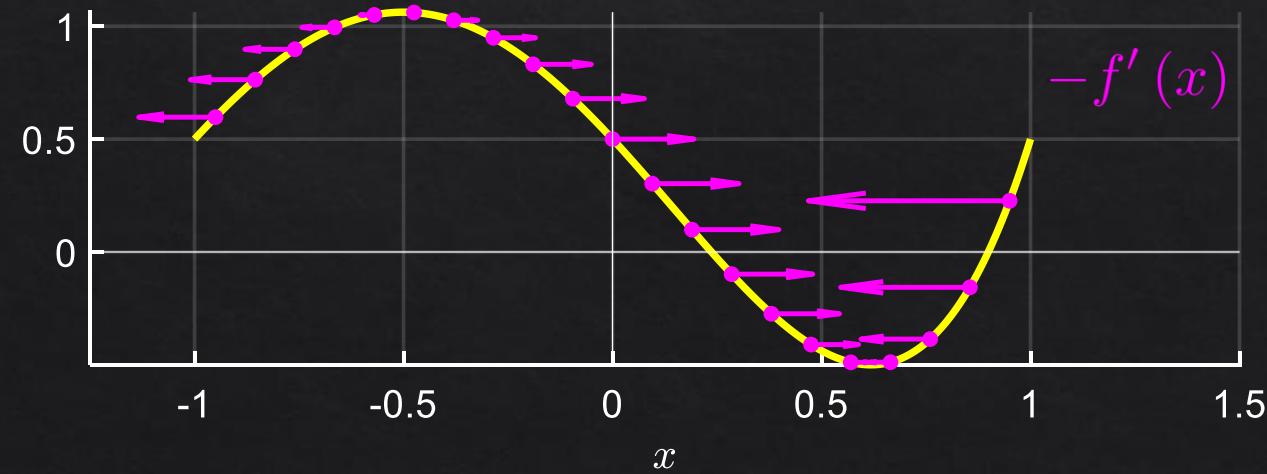
$$f(x) = x^4 + 2x^3 - x^2 - 2x + \frac{1}{2}$$



short arrow $\implies f'(x) \approx 0$ $\implies f$ is flat (around x)

long arrow $\implies |f'(x)|$ is large \implies steep slope (around x)

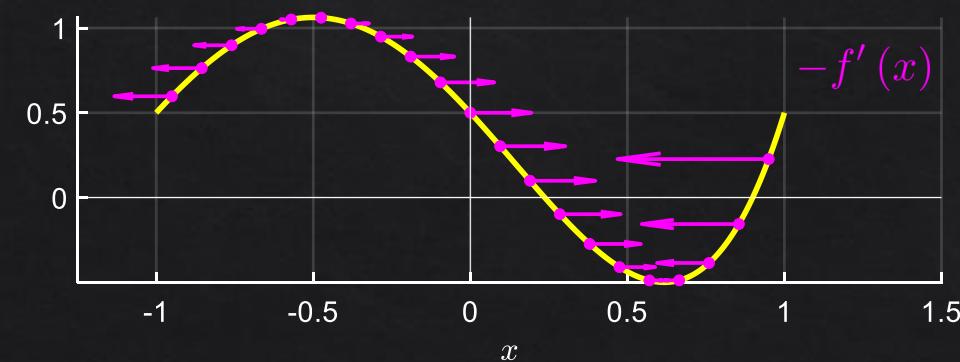
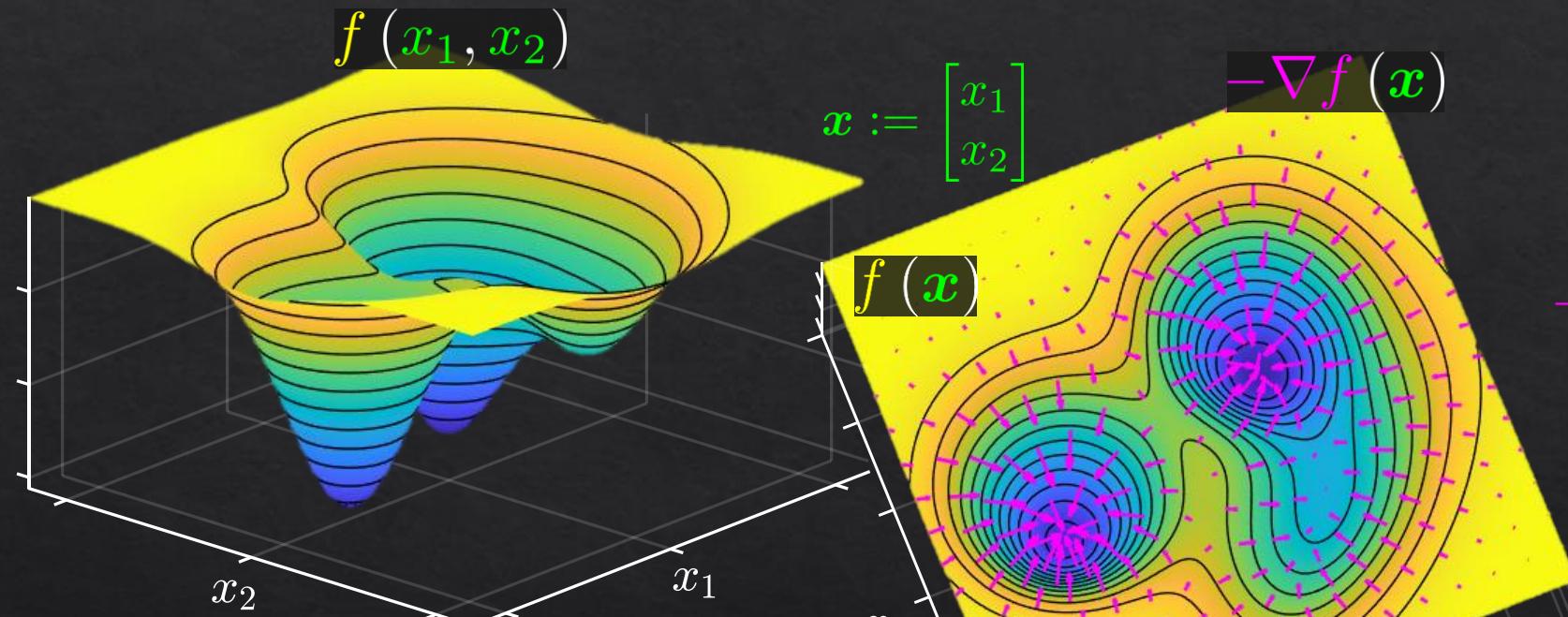
- The arrows indicate the size of the derivative $f'(x)$ (including sign).
- Note that $-f'(x)$ points in a descent direction:



The Gradient

- Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$.
- The gradient $\nabla f(\mathbf{x})$ is a vector of all partial derivatives:

$$\nabla f(\mathbf{x}) := \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \frac{\partial f(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_d} \end{bmatrix} \in \mathbb{R}^d$$



$-\nabla f(\mathbf{x})$ points in a descent direction

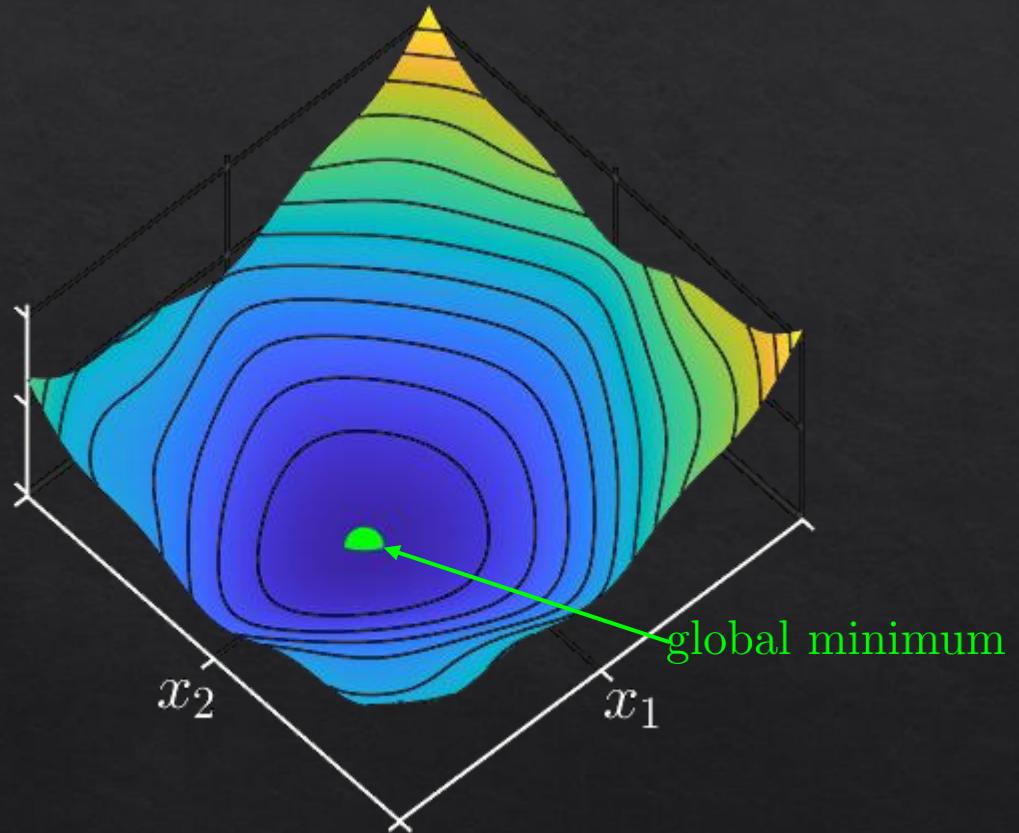
$\|\nabla f(\mathbf{x})\| \approx 0 \implies f$ is flat (around \mathbf{x})

$\|\nabla f(\mathbf{x})\|$ is large \implies steep slope (around \mathbf{x})

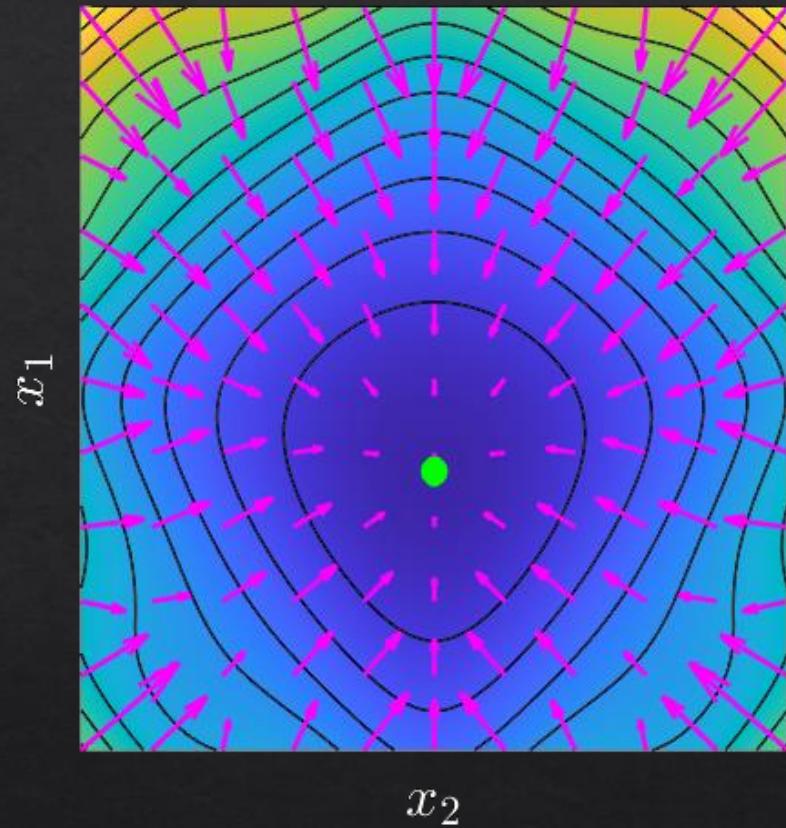
Gradient (Steepest) Descent

- Consider the following function:

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



$$-\nabla f(\mathbf{x})$$



- $-\nabla f$ points in the direction of (locally) maximum decrease in f .
 - We can reach the global minimum by following the (minus) gradient.
- (*) – Note that f is not convex.

Gradient (Steepest) Descent – Algorithm

- Input: A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.
- Output: A local minimum $\mathbf{x}^\star \in \mathbb{R}^d$.

Step I:

Set an initial point $\mathbf{x}_0 \in \mathbb{R}^d$.

Step II:

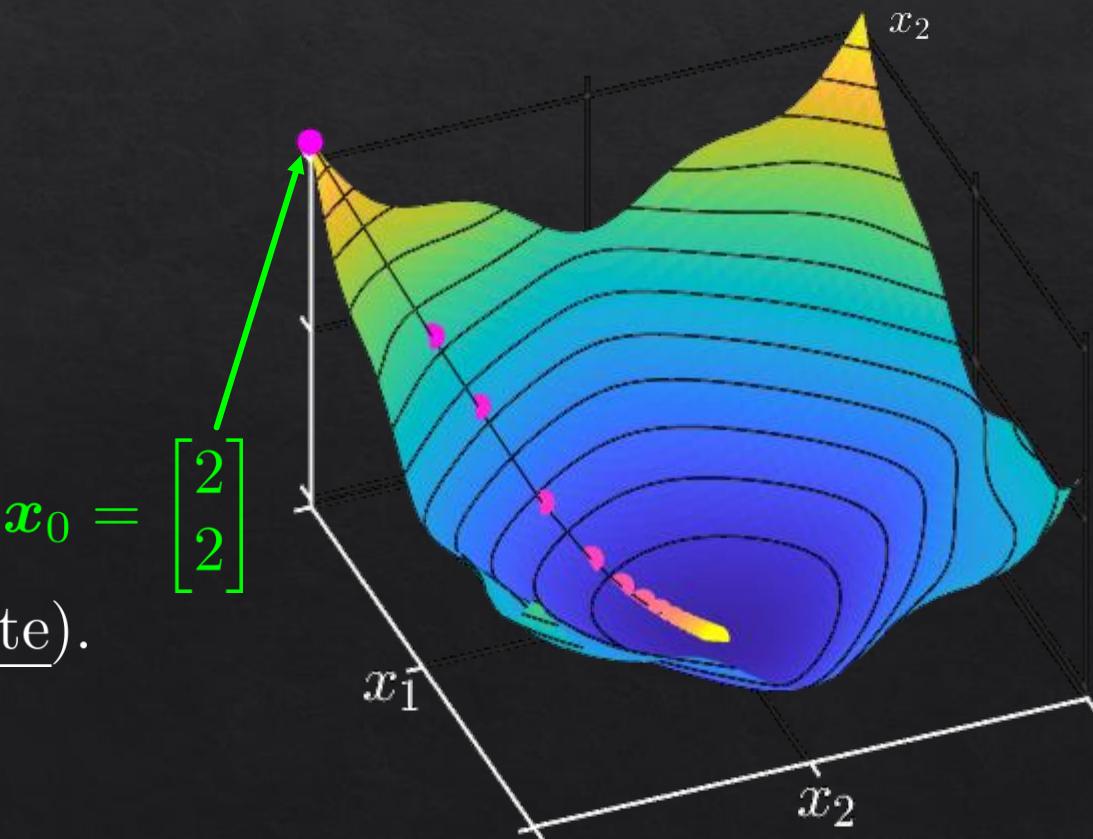
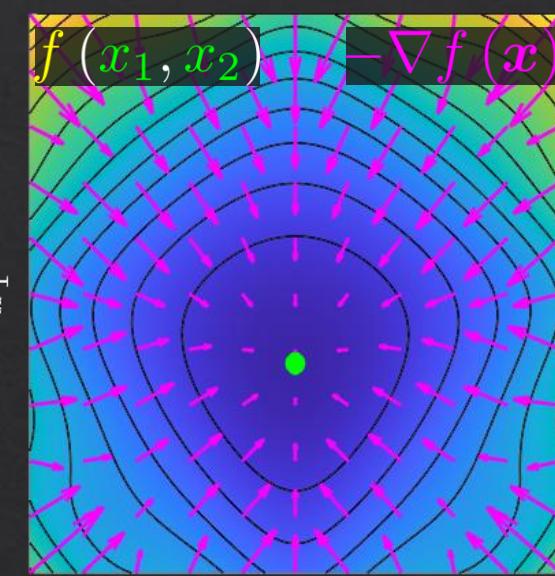
1. **for** $k = 0, 1, 2, \dots$ Update:

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

2. **return** $\mathbf{x}^\star = \mathbf{x}_k$.

- $\mu \in \mathbb{R}$ is the step size (a.k.a. the learning rate).

$$\mathbf{x}_0 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$



Gradient Descent – Example



Gradient Descent

Step I:

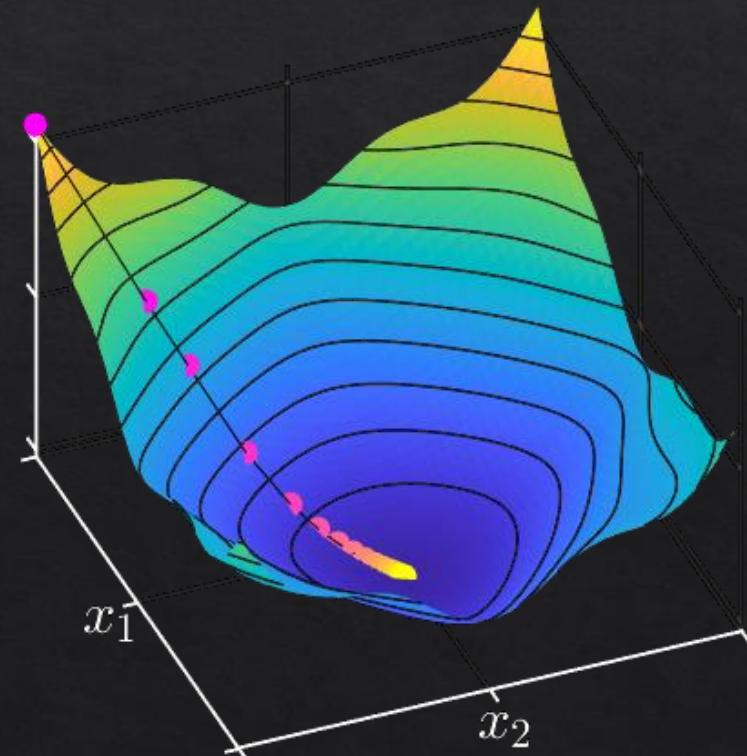
- Set an initial point $\mathbf{x}_0 \in \mathbb{R}^d$.

Step II:

1. **for** $k = 0, 1, 2, \dots$ Update:

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

2. **return** $\mathbf{x}^\star = \mathbf{x}_k$.



Matrix Inversion Example

I – Original Image



$$\mathbf{I} \in \mathbb{R}^{256 \times 256}$$

I_b – Blurred Image



$$\mathbf{I}_b \in \mathbb{R}^{256 \times 256}$$

$$\mathbf{I} = \begin{bmatrix} | & & | \\ \mathbf{i}_1 & \cdots & \mathbf{i}_{256} \\ | & & | \end{bmatrix},$$

column
stack

$$\mathbf{x} := \mathbf{I}^{\text{cs}} =$$

$$\begin{bmatrix} | \\ \mathbf{i}_1 \\ | \\ \vdots \\ | \\ \mathbf{i}_{256} \\ | \end{bmatrix}$$

$$\in \mathbb{R}^{256^2}$$

$$256^2 = 65,536$$

blurring
matrix

$$\mathbf{y} := \mathbf{I}_b^{\text{cs}} \in \mathbb{R}^{256^2}$$

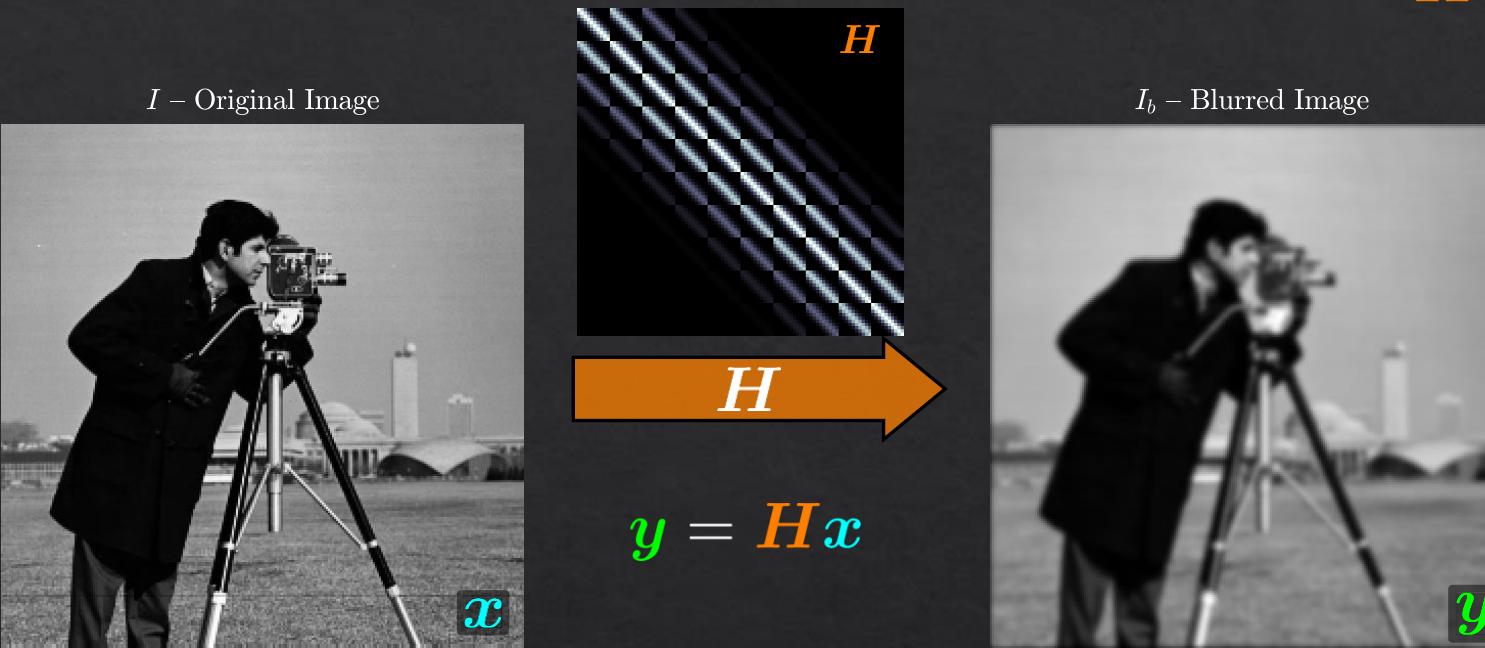
$$\mathbf{y} = \mathbf{H}\mathbf{x}$$

$$\mathbf{H} \in \mathbb{R}^{65,536 \times 65,536}$$

$$\begin{bmatrix} | \\ \mathbf{y} \\ | \end{bmatrix} = \begin{bmatrix} | \\ \mathbf{H} \\ | \end{bmatrix} \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix}$$

Matrix Inversion Example

$$\mathbf{H} \in \mathbb{R}^{65,536 \times 65,536}$$



- Given \mathbf{y} and \mathbf{H} , find \mathbf{x} .
- The solution is: $\mathbf{x} = \mathbf{H}^{-1}\mathbf{y}$
- Since $\mathbf{H} \in \mathbb{R}^{65,536 \times 65,536}$, standard computers cannot compute \mathbf{H}^{-1} .
- Using gradient descent we can solve:

$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2 \quad \implies \quad \mathbf{x}^* = \mathbf{x}$$

Matrix Inversion Example

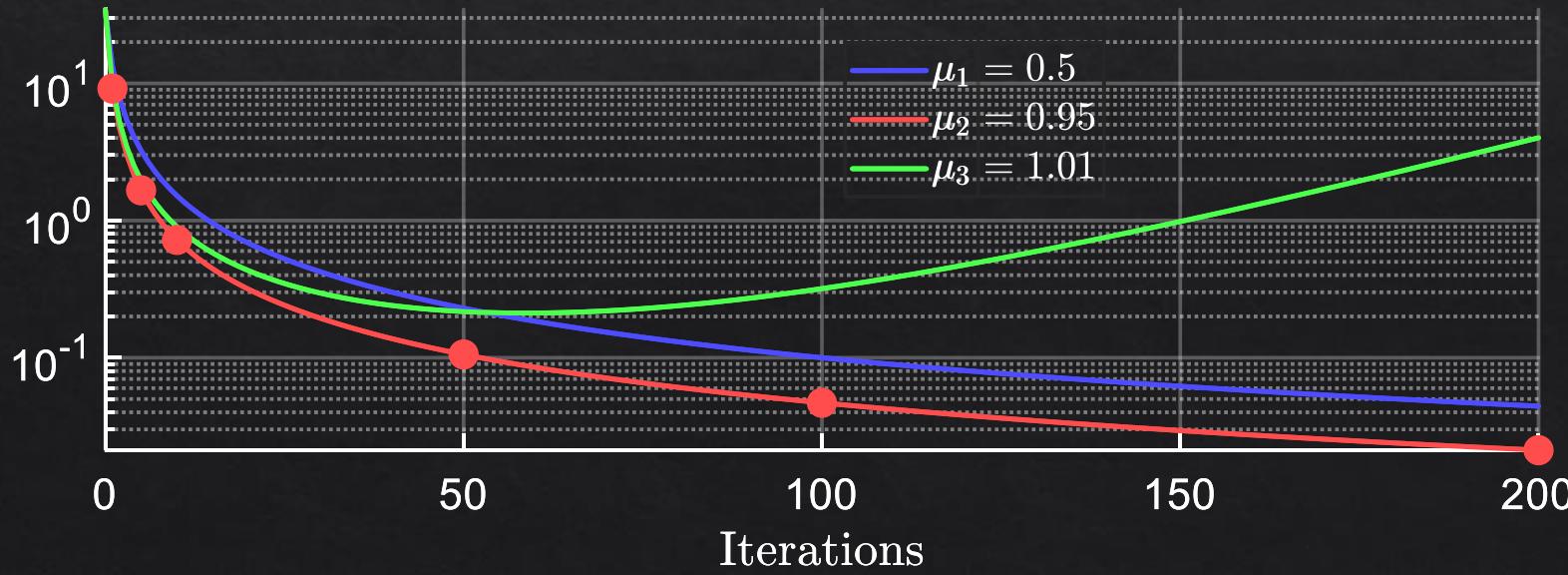
$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \underbrace{\|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2}_{:= f(\tilde{\mathbf{x}})}$$

- The gradient is given by:

$$\left. \begin{array}{l} f(\mathbf{x}) = (\mathbf{y} - \mathbf{H}\mathbf{x})^2 \\ f'(\mathbf{x}) = -2\mathbf{H}(\mathbf{y} - \mathbf{H}\mathbf{x}) \end{array} \right\} \text{1D}$$

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

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



I – Original Image



H



H

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

I_b – Blurred Image



$$\mathbf{y}$$

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

Can you order μ_1, μ_2, μ_3 from small to large?

$$\mu_1 < \mu_2 < \mu_3$$

Matrix Inversion Example

$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \underbrace{\|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2}_{:= f(\tilde{\mathbf{x}})}$$

- The gradient is given by:

$$\left. \begin{array}{l} f(\mathbf{x}) = (\mathbf{y} - \mathbf{H}\mathbf{x})^2 \\ f'(\mathbf{x}) = -2\mathbf{H}(\mathbf{y} - \mathbf{H}\mathbf{x}) \end{array} \right\} \text{1D}$$

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

I – Original Image



\mathbf{H}



\mathbf{H}

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

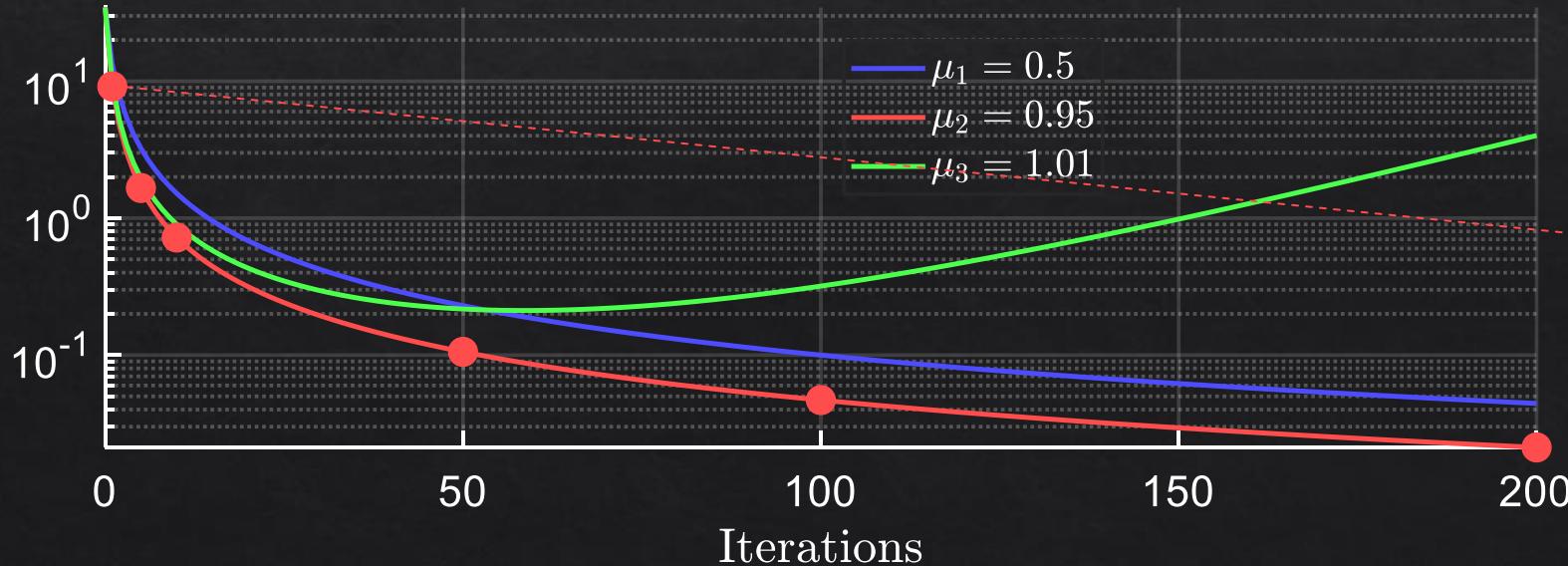
I_b – Blurred Image



\mathbf{y}

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

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



Matrix Inversion Example

$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \underbrace{\|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2}_{:= f(\tilde{\mathbf{x}})}$$

- The gradient is given by:

$$\left. \begin{array}{l} f(\mathbf{x}) = (\mathbf{y} - \mathbf{H}\mathbf{x})^2 \\ f'(\mathbf{x}) = -2\mathbf{H}(\mathbf{y} - \mathbf{H}\mathbf{x}) \end{array} \right\} \text{1D}$$

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

I – Original Image



\mathbf{H}



I_b – Blurred Image

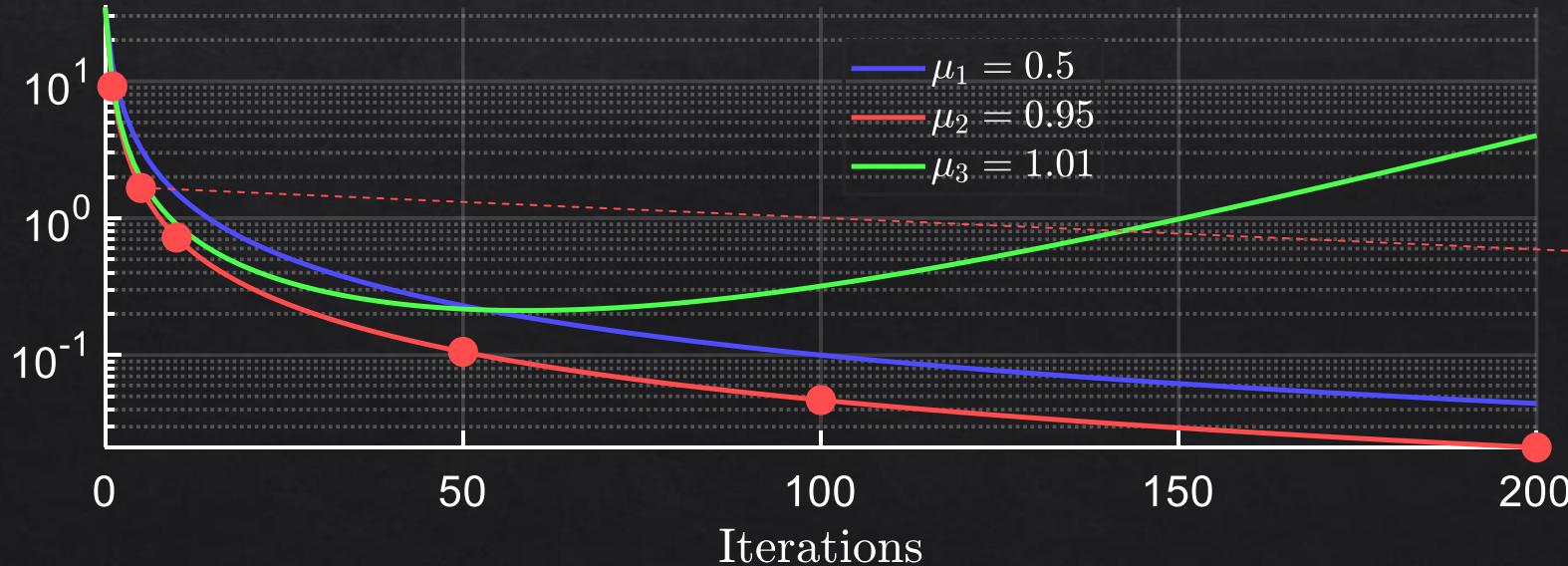
\mathbf{H}



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

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

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



Matrix Inversion Example

$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \underbrace{\|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2}_{:= f(\tilde{\mathbf{x}})}$$

- The gradient is given by:

$$\left. \begin{array}{l} f(\mathbf{x}) = (\mathbf{y} - \mathbf{H}\mathbf{x})^2 \\ f'(\mathbf{x}) = -2\mathbf{H}(\mathbf{y} - \mathbf{H}\mathbf{x}) \end{array} \right\} \text{1D}$$

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

I – Original Image



\mathbf{H}



\mathbf{H}

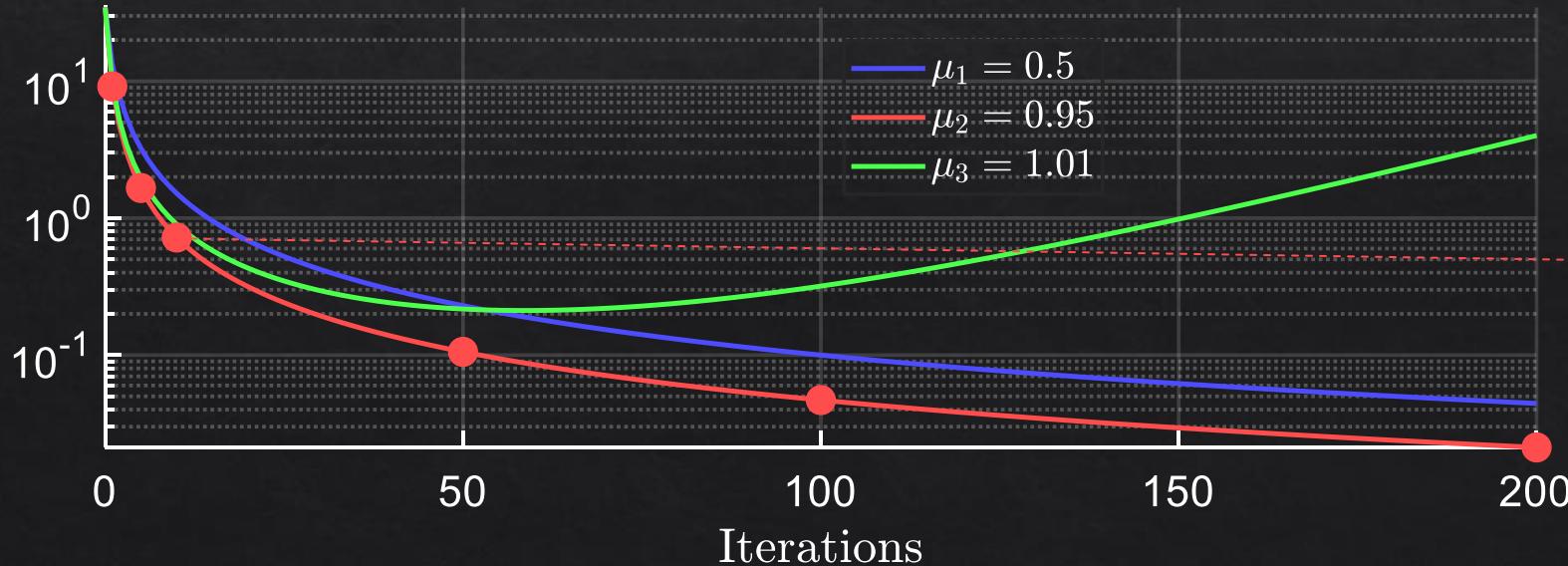
$$\mathbf{y} = \mathbf{H}\mathbf{x}$$

I_b – Blurred Image



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

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



Matrix Inversion Example

$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \underbrace{\|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2}_{:= f(\tilde{\mathbf{x}})}$$

- The gradient is given by:

$$\left. \begin{array}{l} f(\mathbf{x}) = (\mathbf{y} - \mathbf{H}\mathbf{x})^2 \\ f'(\mathbf{x}) = -2\mathbf{H}(\mathbf{y} - \mathbf{H}\mathbf{x}) \end{array} \right\} \text{1D}$$

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

I – Original Image



H



I_b – Blurred Image

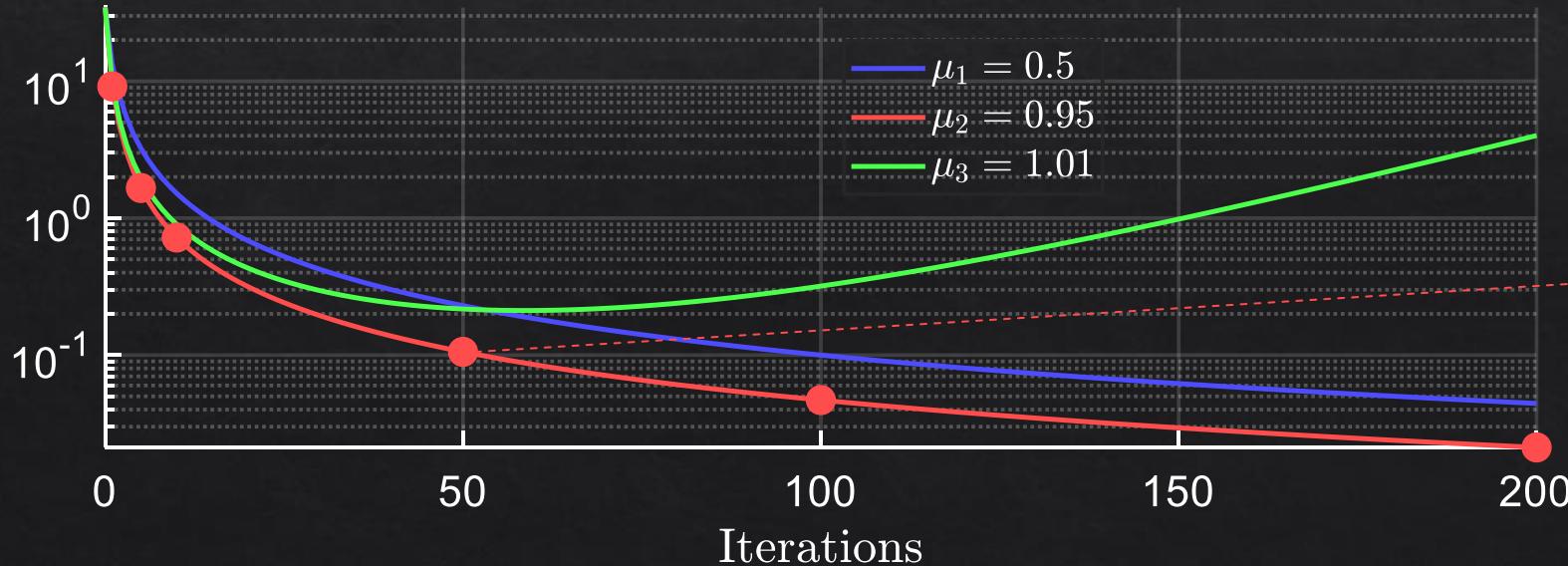
H



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

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

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



Matrix Inversion Example

$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \underbrace{\|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2}_{:= f(\tilde{\mathbf{x}})}$$

- The gradient is given by:

$$\left. \begin{array}{l} f(\mathbf{x}) = (\mathbf{y} - \mathbf{H}\mathbf{x})^2 \\ f'(\mathbf{x}) = -2\mathbf{H}(\mathbf{y} - \mathbf{H}\mathbf{x}) \end{array} \right\} \text{1D}$$

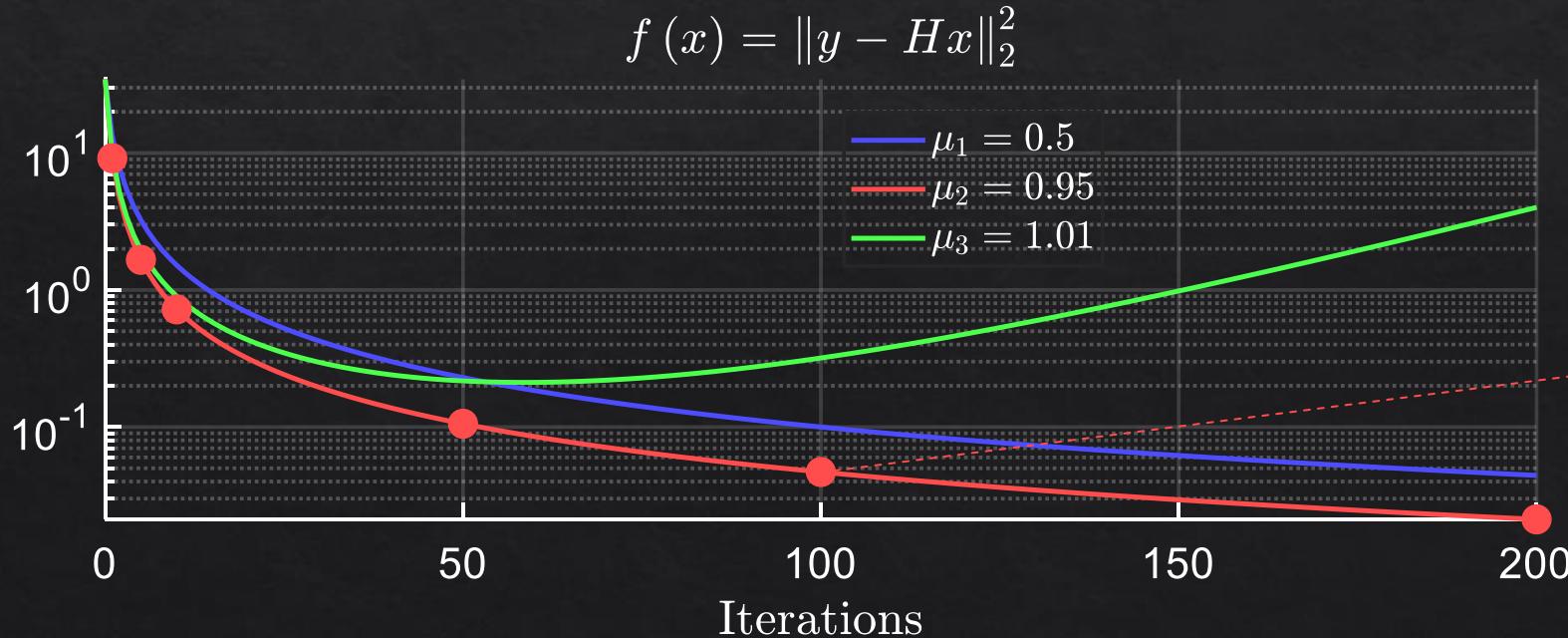
$$\nabla f(\mathbf{x}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{x})$$



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



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



Matrix Inversion Example

$$\mathbf{x}^* = \arg \min_{\tilde{\mathbf{x}}} \underbrace{\|\mathbf{y} - \mathbf{H}\tilde{\mathbf{x}}\|_2^2}_{:= f(\tilde{\mathbf{x}})}$$

- The gradient is given by:

$$\left. \begin{array}{l} f(\mathbf{x}) = (\mathbf{y} - \mathbf{H}\mathbf{x})^2 \\ f'(\mathbf{x}) = -2\mathbf{H}(\mathbf{y} - \mathbf{H}\mathbf{x}) \end{array} \right\} \text{1D}$$

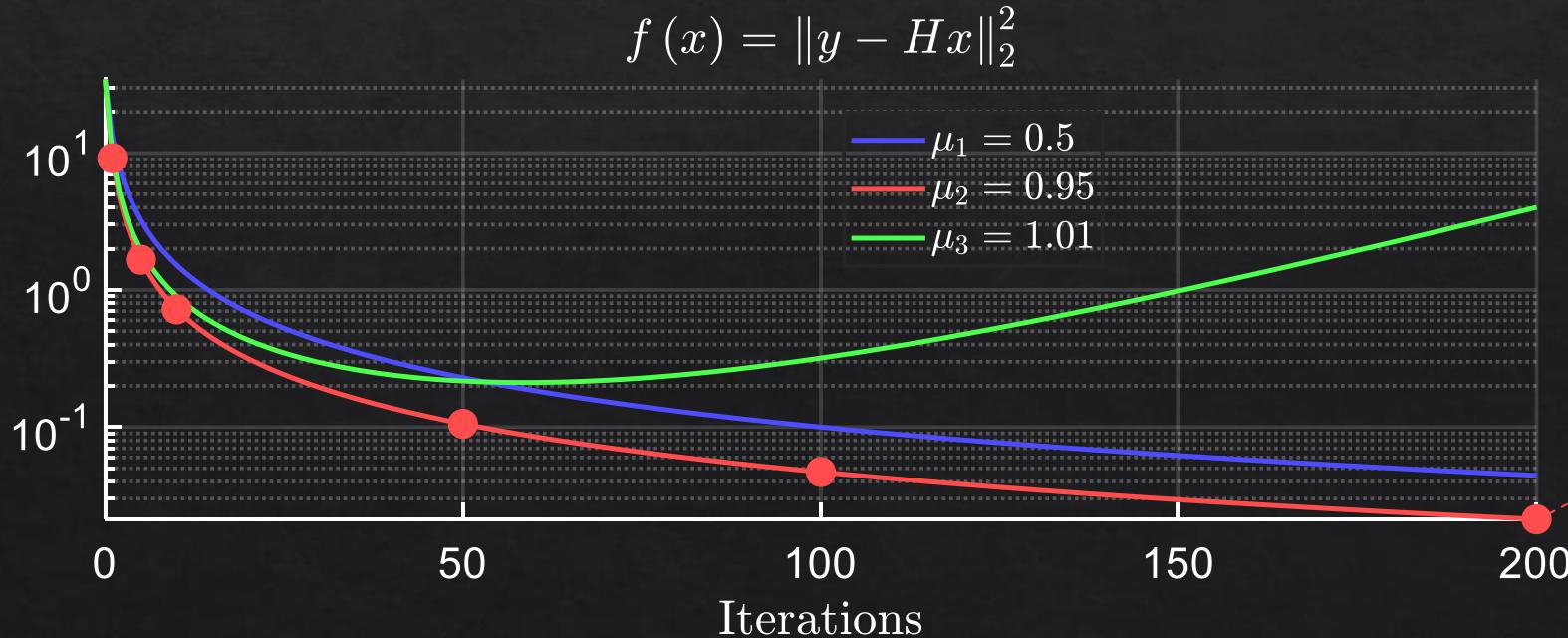
$$\nabla f(\mathbf{x}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{x})$$



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



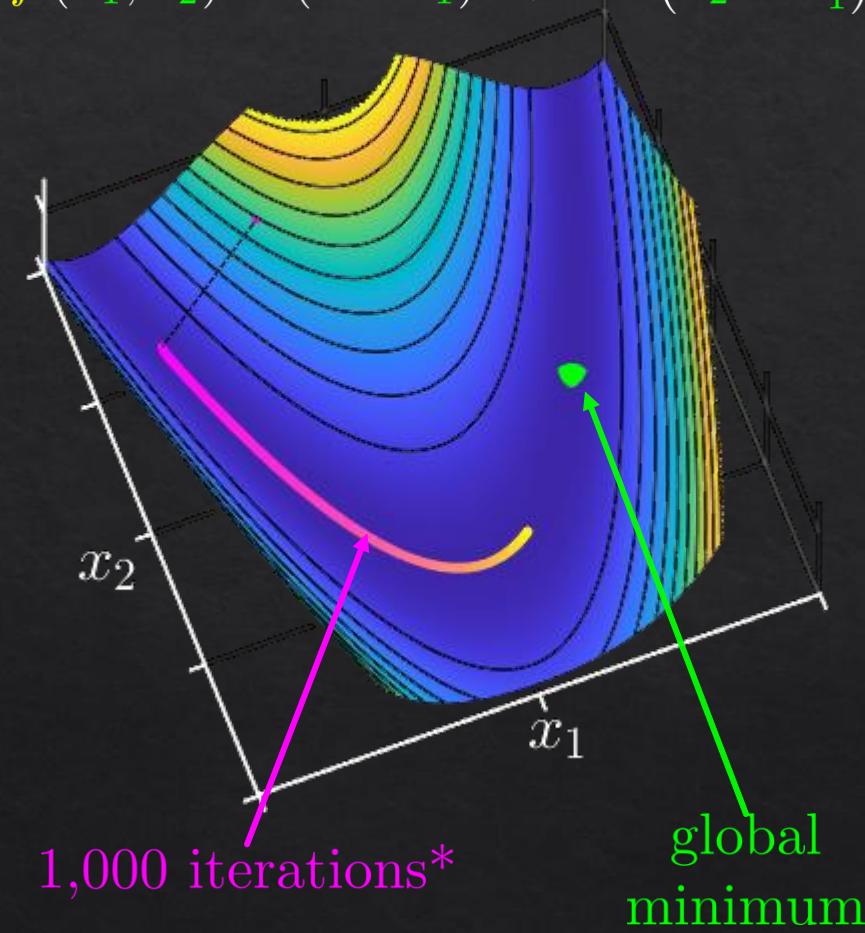
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mu \nabla f(\mathbf{x}_k)$$



Non-convex Functions

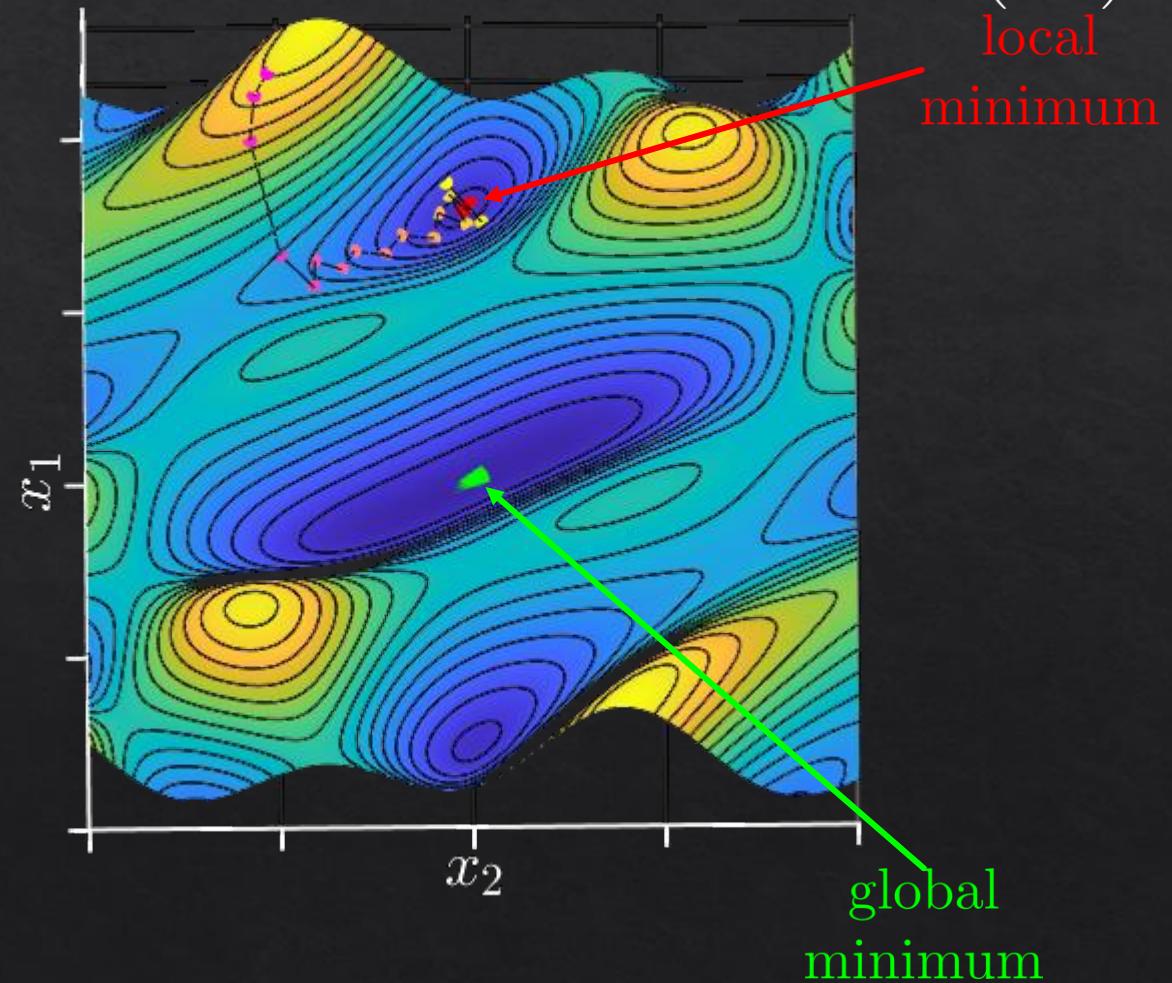
- Slow convergence – Rosenbrock function:

$$f(x_1, x_2) = (1 - x_1)^2 + 100(x_2 - x_1^2)^2$$



- Local minimum:

$$f(x_1, x_2) = \sin^2(x_1 x_2) + \sin^2(2x_1 + x_2) + \left(\frac{1}{5}x_1\right)^2$$



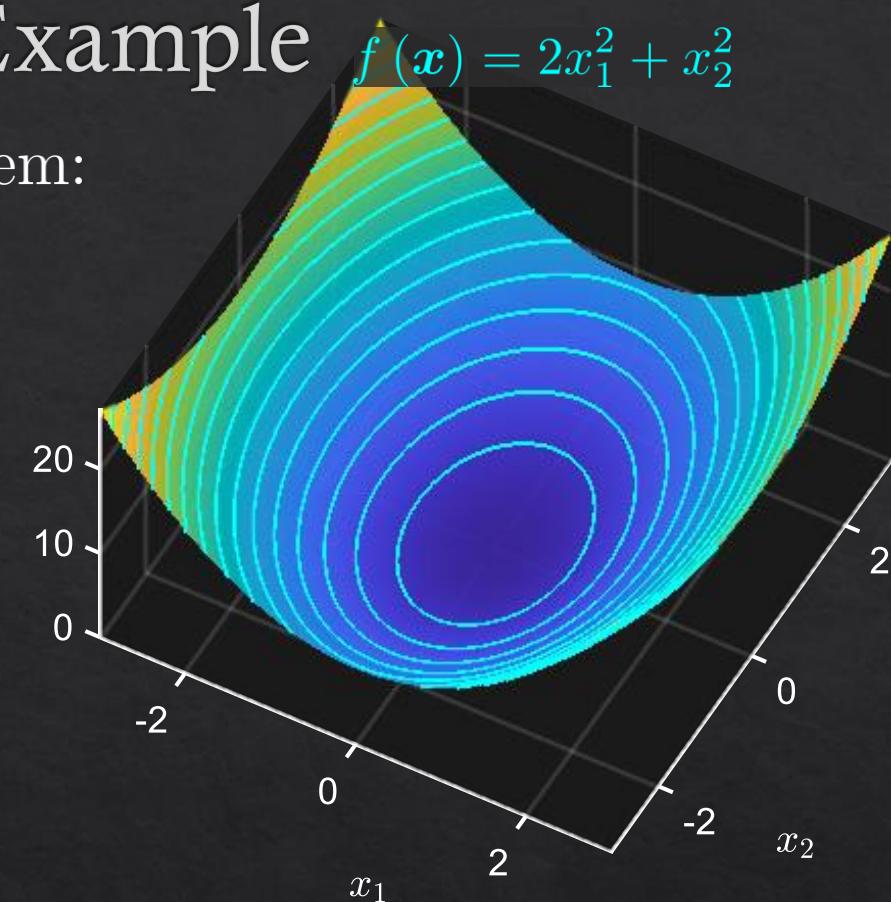
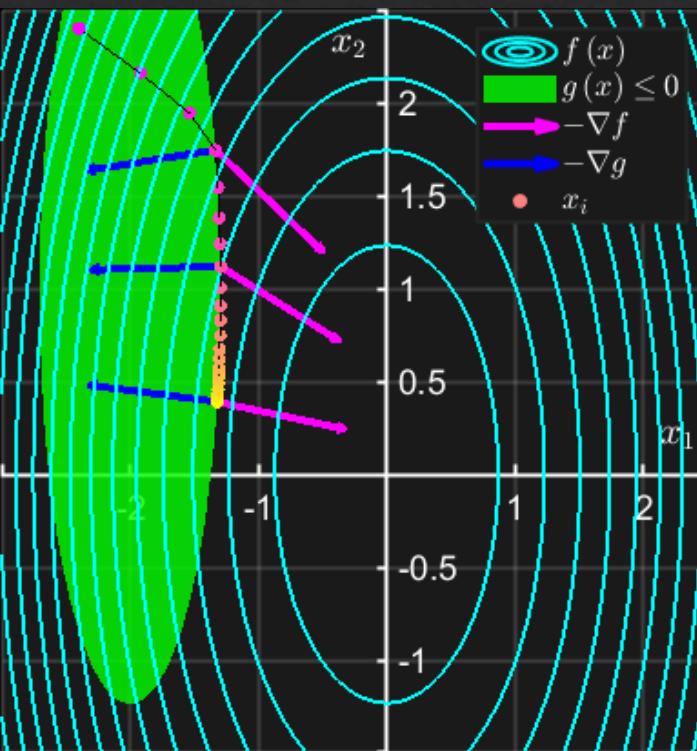
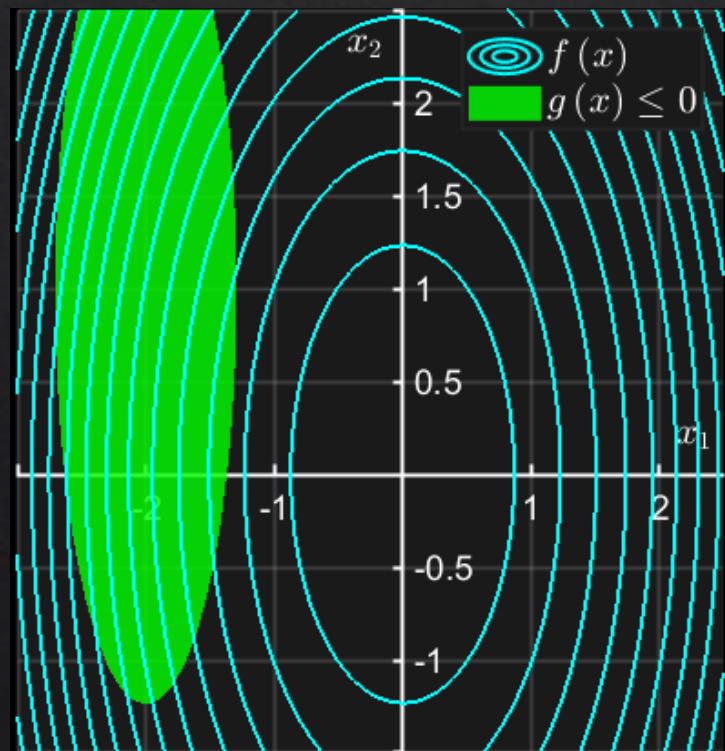
(*) – with the optimal step size.

Constrained Optimization – A Simple Example

$$f(\mathbf{x}) = 2x_1^2 + x_2^2$$

- Consider the following constrained optimization problem:

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



we stop when

$$\nabla f(\mathbf{x}) = -\lambda \nabla g(\mathbf{x})$$

for some $\lambda \geq 0$

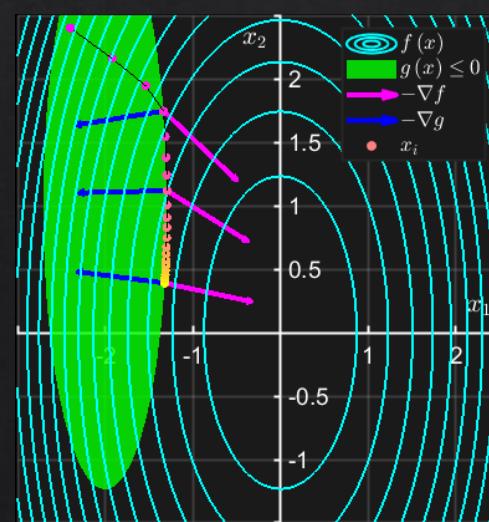
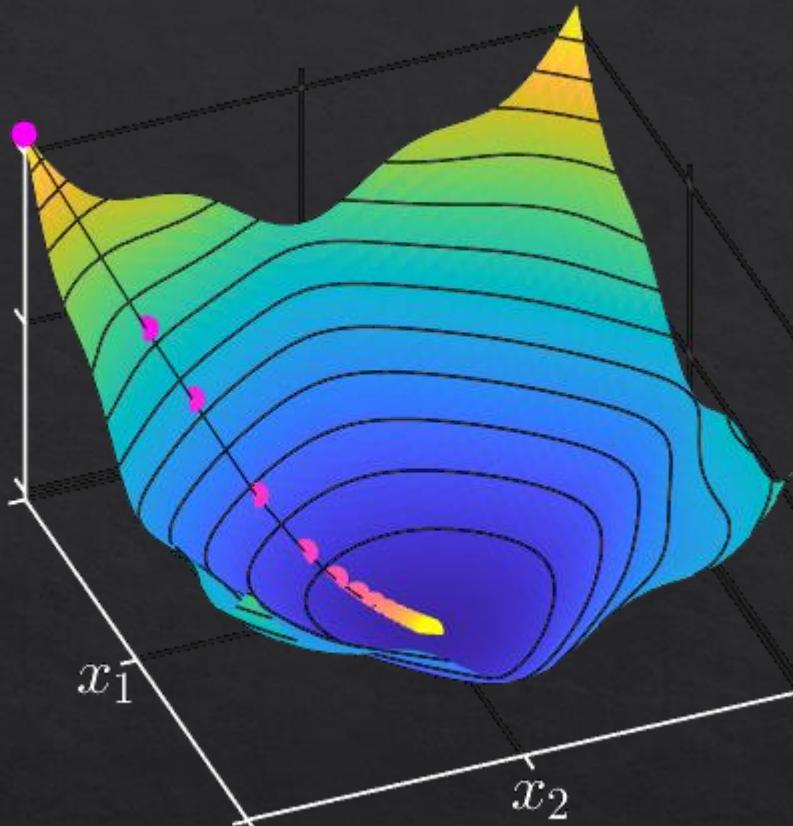
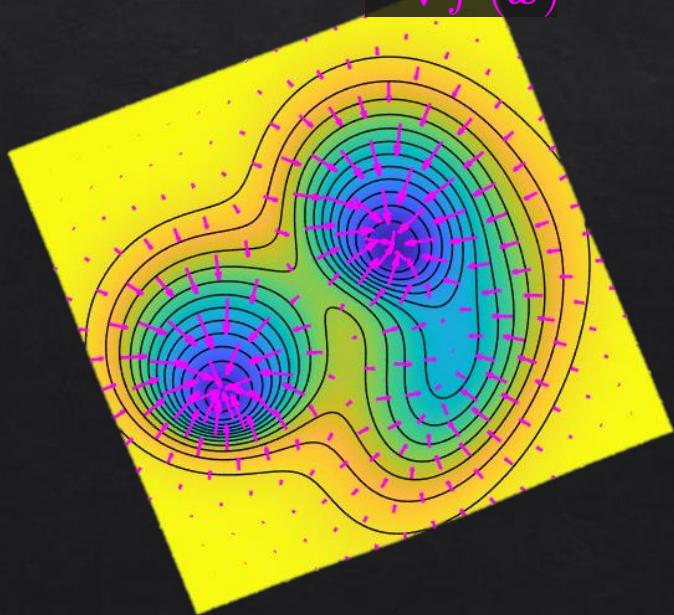
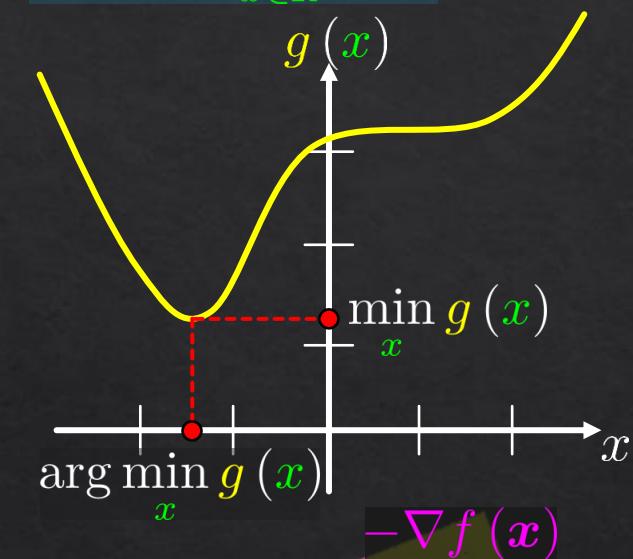
There are several algorithms for such problems (such as the augmented Lagrangian methods).

Other Important Topics

- Hessian and condition number
- Adaptive step size (and momentum)
- Line search
- Conjugate gradient
- Stochastic gradient descent
- Newton method (and quassi Newton)
- minimax problems
- The Lagrangian
- Dual problem

Questions

$$x^* = \arg \min_{x \in \mathbb{R}} f(x)$$



fixelalgorithms.gitlab.io