

# Machine Learning

Gradient Descent, Newton's Method

# System of Linear Equations – Infeasible

Consider  $\mathbf{A} \in \mathbb{C}^{M \times N}$ . We want to **solve** the system of equations ( $M$  equations,  $N$  unknowns):

$$\mathbf{y} = \mathbf{A}\mathbf{x}.$$

That is, we want to find  $\mathbf{x}$  such that the above equation is satisfied, given  $\mathbf{A}$  and  $\mathbf{y}$ .

Case 1: Infeasible (no solution). The system is solvable iff  $\mathbf{y} \in \text{span}(\mathbf{A})$ . A sufficient (but not necessary) condition is  $\text{rank}(\mathbf{A}) = M$ .

# SLE – General SVD Solution

Case 2: Feasible. If  $\mathbf{y} \in \text{span}(\mathbf{A}) = \text{span}(\mathbf{U})$ , the system  $\mathbf{y} = \mathbf{A}\mathbf{x}$  is solvable.

For RSVD  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H$ , the general solution is given by

$$\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y} + \mathbf{b}$$

for any  $\mathbf{b} \in \mathcal{N}(\mathbf{A}) = \mathcal{N}(\mathbf{V}^H) = \text{span}(\mathbf{V})^\perp = \{\mathbf{z} \in \mathbb{C}^N : \mathbf{z}^H\mathbf{y} = 0 \ \forall \mathbf{y} \in \text{span}(\mathbf{V})\}$ .

Verify:  $\mathbf{A}\mathbf{x} = \mathbf{A}(\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y} + \mathbf{b}) = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^H(\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y} + \mathbf{b}) = \mathbf{U}\mathbf{U}^H\mathbf{y} = \mathbf{y}$ .

Thus, there are as many solutions as  $|\mathcal{N}(\mathbf{A})|$ . Subspaces come in two cardinalities: 0 and infinity.

# SLE – Full Row-Rank

Special case 1:  $\rho = \text{rank}(\mathbf{A}) = M < N$  (full-row rank). This implies that  $\mathbf{A}$  is wide. Therefore, there are  $|\mathcal{N}(\mathbf{A})| = \infty$  solutions to the system. Moreover,

$$\begin{aligned}\mathbf{A}^{\dagger R} &= \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1} \\ &= \mathbf{V} \mathbf{\Sigma} \mathbf{U}^H (\mathbf{U} \mathbf{\Sigma} \mathbf{V}^H \mathbf{V} \mathbf{\Sigma} \mathbf{U}^H)^{-1} \\ &= \mathbf{V} \mathbf{\Sigma} \mathbf{U}^H (\mathbf{U}^H)^{-1} \mathbf{\Sigma}^{-2} (\mathbf{U})^{-1} \\ &= \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^H.\end{aligned}$$

More unknowns ( $N$ ) than independent equations ( $\rho = M$ ).

In this special case, the system solution takes the special form:

$$\mathbf{x} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^H \mathbf{y} + \mathbf{b} = \mathbf{A}^{\dagger R} \mathbf{y} + \mathbf{b}, \quad \forall \mathbf{b} \in \mathcal{N}(\mathbf{A})$$

and no decomposition of  $\mathbf{A}$  (e.g., SVD) is needed for solving the system.

# SLE – Full Column-Rank

Special case 2:  $\rho = \text{rank}(\mathbf{A}) = N < M$  (full-column rank). This implies that  $\mathcal{N}(\mathbf{A}) = \mathbf{0}_N$  and the system has a single solution. Moreover,

$$\begin{aligned}\mathbf{A}^{\dagger L} &= (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \\ &= (\mathbf{V} \Sigma \mathbf{U}^H \mathbf{U} \Sigma \mathbf{V}^H)^{-1} \mathbf{V} \Sigma \mathbf{U}^H \\ &= (\mathbf{V}^H)^{-1} \Sigma^{-2} (\mathbf{V})^{-1} \mathbf{V} \Sigma \mathbf{U}^H \\ &= \mathbf{V} \Sigma^{-1} \mathbf{U}^H.\end{aligned}$$

More independent equations to satisfy ( $M$ ) than unknowns to tune ( $N$ ).

In this special case,

$$\mathbf{x} = \mathbf{V} \Sigma^{-1} \mathbf{U}^H \mathbf{y} + \mathbf{0}_N = \mathbf{A}^{\dagger L} \mathbf{y}$$

and no decomposition of  $\mathbf{A}$  (e.g., SVD) is needed for solving the system.

## SLE – Square, Full Rank

Special case 3:  $\rho = \text{rank}(\mathbf{A}) = M = N$  (square full rank). This implies that  $\mathcal{N}(\mathbf{A}) = \mathbf{0}_N$  and the system again has a single solution. Moreover,

$$\begin{aligned}\mathbf{A}^{\dagger L} &= \mathbf{A}^{\dagger L} = \mathbf{A}^{-1} \\ &= (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^H)^{-1} \\ &= (\mathbf{V}^H)^{-1}\mathbf{\Sigma}^{-1}(\mathbf{U})^{-1} \\ &= \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\end{aligned}$$

In this special case,

$$\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y} + \mathbf{0}_N = \mathbf{A}^{-1}\mathbf{y}$$

and, again, no decomposition of  $\mathbf{A}$  (e.g., SVD) is needed for solving the system.

# Multivariate Derivative

Simple extension of the concept of univariate differentiation.

$$\frac{\partial X_{k,l}}{\partial X_{i,j}} = \delta_{i,k} \delta_{l,j}$$

Vector derivative w.r.t. scalar  $\rightarrow$  vector  $\mathbf{a} = \frac{\partial x}{\partial \mathbf{y}}$ ,  $a_i = \frac{\partial x_i}{\partial y}$

Scalar derivative w.r.t. vector  $\rightarrow$  vector  $\mathbf{a} = \frac{\partial x}{\partial \mathbf{y}}$ ,  $a_i = \frac{\partial x}{\partial y_i}$ . Also written as  $\mathbf{a} = \nabla_{\mathbf{y}} x$ .

Vector derivative w.r.t. vector  $\rightarrow$  matrix  $\mathbf{A} = \frac{\partial \mathbf{x}}{\partial \mathbf{y}}$ ,  $A_{i,j} = \frac{\partial x_i}{\partial y_j}$

Extends to matrix w.r.t to martrix...

(see more here: <https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>)

# Multivariate Derivative - Examples

Example 1: First order (Linear)

$$\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{y} = \mathbf{a}$$

$$a_i = \frac{\partial(\mathbf{x}^T \mathbf{y})}{\partial x_i} = \frac{\partial(\sum_j x_j y_j)}{\partial x_i} = y_i \Rightarrow \mathbf{a} = \mathbf{y}$$

Example 2: Second order (Quadratic)

$$\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{a}$$

$$a_i = \frac{\partial(\mathbf{x}^T \mathbf{A} \mathbf{x})}{\partial x_i} = \frac{\partial(\sum_j \sum_k x_j A_{j,k} x_k)}{\partial x_i} = \dots = \sum_j (A_{i,j} + A_{j,i}) x_j \Rightarrow \mathbf{a} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$



# Multivariate Optimization

$$P: \min_{x \in F} f(x)$$

This is a multivariate optimization problem.

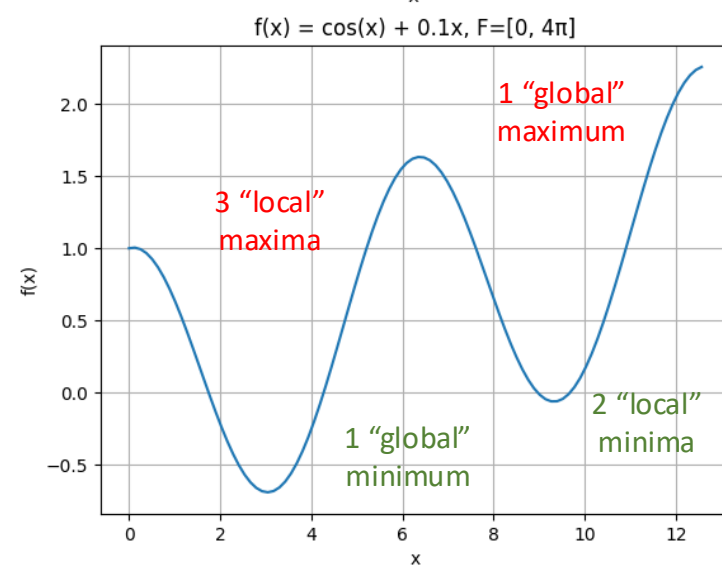
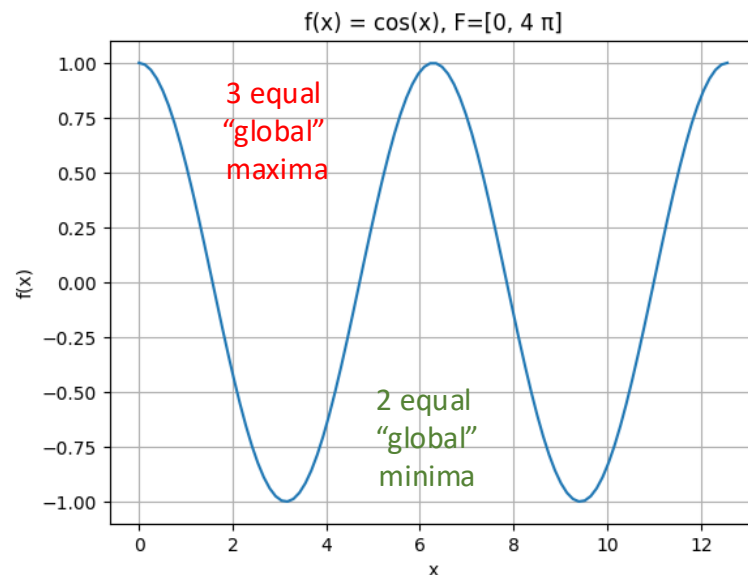
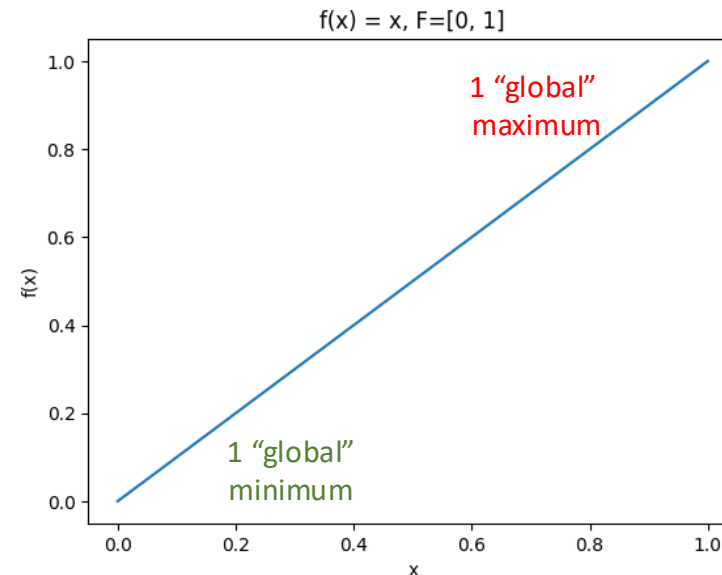
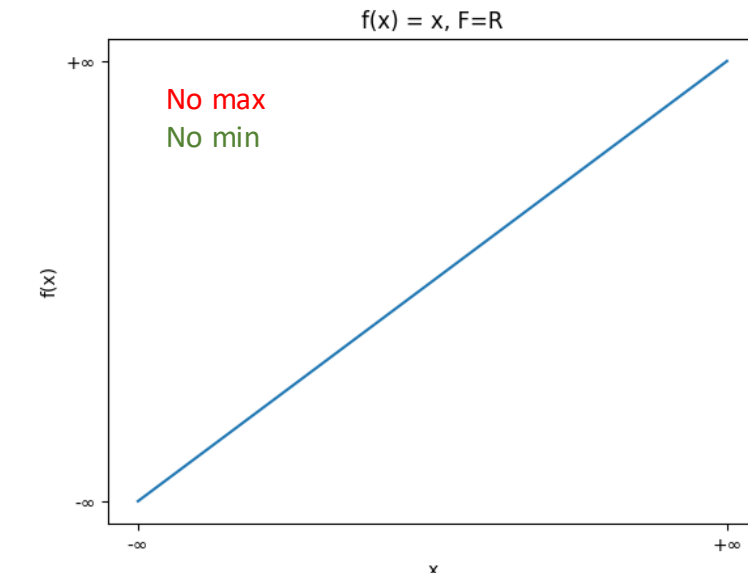
Find the argument  $x$  that minimizes objective function  $f: F \rightarrow \mathbb{R}$

Can have no solution, unique solution, multiple solutions, infinite solutions.

Examples:

- $f(x) = x, F = \mathbb{R}$ . No solution.
- $f(x) = x, F = [0,1]$ . Single solution.
- $f(x) = x^2, F = \mathbb{R}$ . Single solution.
- $f(x) = \cos(x), F = [0, 4\pi]$ . Two solutions.
- $f(x) = \cos(x), F = \mathbb{R}$ . Infinitely many solutions.

# Multivariate Optimization



# Multivariate Optimization

$$P: \min_{x \in F} f(x)$$

Therefore, the combination of  $F$  and  $f$  determine whether a problem is solvable or not; if it is solvable, they determine how many solutions it has and how easy/difficult it is to obtain any of the solutions.

Convex  $F$  and  $f$  make  $P$  an instance of **convex optimization**.

This is an important field of study.

**Convex optimization problems have a single global minimum. No local minima.**

See more here: <https://web.stanford.edu/~boyd/cvxbook/>

# Least Squares (LS)

Least Squares problem: Given  $A$  and  $y$ , solve:

$$P: \min_{x \in \mathbb{R}^N} \|y - Ax\|_2^2$$

This is an instance of **convex** optimization (quadratic minimization objective, unconstrained).

Need an algorithm (i.e., methodology) for solving this problem (i.e., find min and argmin).

We start with inspecting the problem, in search of an efficient solution.

Based on our SLE understanding, there is a case where the algorithm can be fast: single-shot solution.

This is the case of  $y \in \text{span}(A)$ ...

## LS – Special Case: $\mathbf{y} \in \text{span}(\mathbf{A})$

We first notice that:  $\|\mathbf{y} - \mathbf{Ax}\|_2^2 \geq 0$

Then, we recall that when  $\mathbf{y} \in \text{span}(\mathbf{A})$ , there exists  $\mathbf{x}$  so that  $\mathbf{y} = \mathbf{Ax}$ .

Specifically, by SLE theory (considering R-SVD  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ ) it holds:

$$\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y} + \mathbf{b}, \mathbf{b} \in \text{null}(\mathbf{A})$$

$$\Rightarrow \mathbf{Ax} = \mathbf{y}$$

$$\Leftrightarrow \mathbf{Ax} - \mathbf{y} = \mathbf{0}$$

$$\Leftrightarrow \|\mathbf{y} - \mathbf{Ax}\|_2^2 = 0$$

## LS – Special Case: $\mathbf{y} \in \text{span}(\mathbf{A})$

Therefore, when  $\mathbf{y} \in \text{span}(\mathbf{A})$ , for any  $\mathbf{b} \in \text{null}(\mathbf{A})$ ,

$$\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y} + \mathbf{b} \in \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$$

And

$$\min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = 0$$

This is an instance of **convex** optimization (quadratic minimization objective, unconstrained feasibility).

1) Perform R-SVD:  $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}) \xleftarrow{RSVD} \mathbf{A}$

2) Return  $\mathbf{x}_{\text{opt}} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y} + \mathbf{b}$  for any  $\mathbf{b} \in \mathcal{N}(\mathbf{A})$  // Or just return  $\mathbf{x}_{\text{opt}} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y}$

- Note: Clearly, the LS problem in this case has infinitely many optimal arguments since

$|\mathcal{N}(\mathbf{A})| = \infty$ . Since  $\mathbf{0} \in \text{null}(\mathbf{A})$ ,  $\mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^H\mathbf{y}$  is one of the infinitely many optimal arguments.

# LS – General Case

Given  $A$  and  $y \in \mathbb{R}^D$ , solve:

$$P: \min_{x \in \mathbb{R}^N} \|y - Ax\|_2^2$$

This is an instance of **convex** optimization (quadratic minimization objective, unconstrained feasibility).

In calculus you learned that the stationary point of a function is where the function derivative becomes 0.

For unconstrained convex optimization (like the LS above), this stationary point is an argmin.

## LS – General Case

So, in the general case, LS can be solved by setting derivative of objective  $L(\mathbf{x}) = \|\mathbf{y} - \mathbf{Ax}\|_2^2$  to 0.

To that end, we expand the quadratic objective as follows:

$$L(\mathbf{x}) = \|\mathbf{y} - \mathbf{Ax}\|_2^2 = (\mathbf{y} - \mathbf{Ax})^T (\mathbf{y} - \mathbf{Ax}) = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - \mathbf{y}^T \mathbf{Ax} - \mathbf{x}^T \mathbf{A}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$$

We calculate the derivative (gradient) as:

$$\begin{aligned}\nabla_{\mathbf{x}} L(\mathbf{x}) &= \nabla_{\mathbf{x}} (\mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - \mathbf{y}^T \mathbf{Ax} - \mathbf{x}^T \mathbf{A}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \\ &= \nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} - \nabla_{\mathbf{x}} \mathbf{y}^T \mathbf{Ax} - \nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A}^T \mathbf{y} + \nabla_{\mathbf{x}} \mathbf{y}^T \mathbf{y} \\ &= (\mathbf{A}^T \mathbf{Ax} - \mathbf{A}^T \mathbf{y})2\end{aligned}$$



## LS – General Case

$$\nabla_{\mathbf{x}} L(\mathbf{x}) = (\mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{A}^T \mathbf{y})2$$

Setting the gradient equal to 0:

$$\nabla_{\mathbf{x}} L(\mathbf{x}) \Big|_{\mathbf{x}_{opt}} = 0$$

$$(\mathbf{A}^T \mathbf{A} \mathbf{x}_{opt} - \mathbf{A}^T \mathbf{y})2 = 0$$

$$\mathbf{A}^T \mathbf{y} = \mathbf{A}^T \mathbf{A} \mathbf{x}_{opt}$$

This is an SLE. Since we do not want to make assumptions on the rank of  $\mathbf{A}$  and  $\mathbf{A}^T \mathbf{A}$ , the general solution will be derived by the R-SVD of  $\mathbf{A}^T \mathbf{A}$ .

## LS – General Case

We start with the R-SVD  $A = U\Sigma V^T$ .

Then we find  $A^T A = (U\Sigma V^T)^T U\Sigma V^T = V\Sigma^T U^T U\Sigma V^T = V\Sigma^2 V^T$ .

This constitutes R-SVD of  $A^T A$ .

Note: R-SVD solution is unique. Once a matrix is written as **orthonormal-columns**  $\times$  **positive diagonal**  $\times$  **orthonormal-rows**, this is the unique R-SVD of that matrix.

We also find that  $\text{null}(A) = \mathcal{N}(A^T A) = \mathcal{N}(V^T)$

Therefore, the solution of this SLE  $A^T y = A^T A x_{opt}$  is:

$$x_{opt} = V\Sigma^{-2}V^H A^T y + b = V\Sigma^{-2}V^T V\Sigma U^T y + b = V\Sigma^{-1}U^T y + b$$

where  $b \in \mathcal{N}(A)$ .

## LS – General Case

So, the solution to the Least Squares problem  $\min_{\mathbf{x} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$  is always

$$\mathbf{x}_{opt} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{y} + \mathbf{b}, \text{ for any } \mathbf{b} \in \mathcal{N}(\mathbf{A})$$

whether  $\mathbf{y}$  is in  $\text{span}(\mathbf{A})$  or not!

If  $\mathbf{y} \in \text{span}(\mathbf{A})$ :  $\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{y} + \mathbf{b}$  minimizes  $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ , setting it equal to 0, since  $\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{y} + \mathbf{b}$  is the solution to the SLE  $\mathbf{y} = \mathbf{A}\mathbf{x}$ .

If  $\mathbf{y} \notin \text{span}(\mathbf{A})$ :  $\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{y} + \mathbf{b}$  minimizes  $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ , which will be positive (not zero), since  $\mathbf{y} = \mathbf{A}\mathbf{x}$  does not have a solution.

So, LS is broader than SLE and boils down to SLE only when  $\mathbf{y} \in \text{span}(\mathbf{A})$

## LS – Special Case (full-row or full-column rank)

When  $\rho = \text{rank}(\mathbf{A}) = M < N$  (full-row rank),  $|\mathcal{N}(\mathbf{A})| = \infty$  and

$$\mathbf{A}^{\dagger R} = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^H.$$

Thus, the LS solution is  $\mathbf{x}_{opt} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^H \mathbf{y} = \mathbf{A}^{\dagger R} \mathbf{y} + \mathbf{b}$ , for any  $\mathbf{b} \in \mathcal{N}(\mathbf{A})$ .

Recall that  $\mathbf{0} \in \mathcal{N}(\mathbf{A})$ , so  $\mathbf{x}_{opt} = \mathbf{A}^{\dagger R} \mathbf{y}$  is one of the infinitely many solutions.

When  $\rho = \text{rank}(\mathbf{A}) = N < M$  (full-column rank),  $\mathcal{N}(\mathbf{A}) = \mathbf{0}_N$  and

$$\mathbf{A}^{\dagger L} = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^H.$$

Thus, the unique LS solution is  $\mathbf{x}_{opt} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^H \mathbf{y} = \mathbf{A}^{\dagger L} \mathbf{y}$ .

# Function Definitions

$L(x)$  is **Continuously Differentiable** (smooth or  $C_1$ ) iff (1)  $\nabla L(x)$  exists for all  $x$  in its domain and (2)  $\nabla L(x)$  is a continuous function.

Example:  $L(x) = x^2$ ,  $\nabla L(x) = 2x$  (continuous everywhere)

$L(x)$  is **Twice Continuously Differentiable** ( $C_2$ ) iff (1) It is continuously differentiable, (2) the second derivative  $\nabla^2 L(x)$  exists for all  $x$  in its domain, and (3) the second derivative  $\nabla^2 L(x)$  is a continuous function.

Example:  $L(x) = x^2$ ,  $\nabla^2 L(x) = 2$  (constant, thus continuous)

$L(x)$  is **Lipschitz Continuous** (LC) iff there exists a constant  $C \geq 0$  such that  $|L(x) - L(y)| \leq C |x - y|$  for all  $x, y$  in its domain.

Example:  $L(x) = \max(0, x)$  (ReLU) is LC w/ Lipschitz constant  $C = 1$ .

$L(x)$  is **Convex** if

$$L(\alpha x + (1 - \alpha)y) \leq \alpha L(x) + (1 - \alpha)L(y)$$

for all  $x, y$  in its domain and  $\alpha \in [0, 1]$ .

Example:  $L(x) = x^2$  (the line segment between any two points on the curve lies above the curve)

# Gradient Descent (GD)

Above we presented closed-form solutions for LS via SVD or pseudoinversions. These solutions are possible because of the convexity of LS.

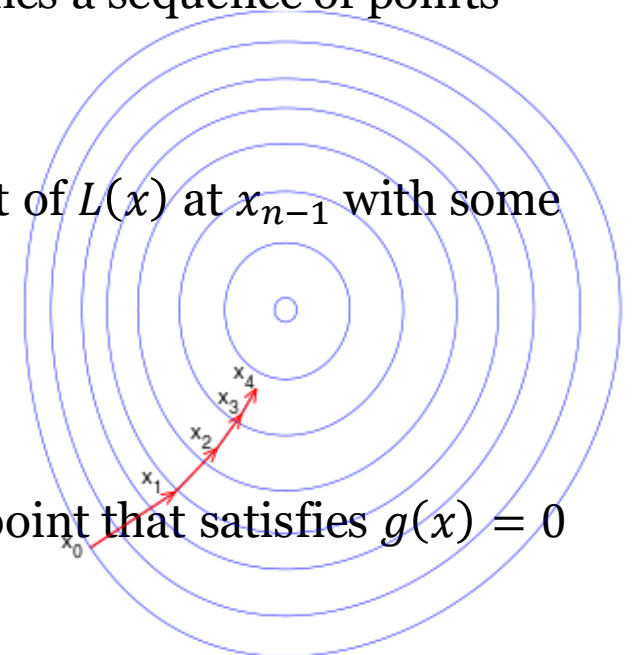
A more general approach to find local/global minima, even for non-convex problems, is the Gradient Descent (GD) method. This is applied when  $L(x)$  is differentiable.

GD is an iterative method that is initialized to some point  $x_0$  and determines a sequence of points  $x_1, x_2, \dots$

When at  $x_{n-1}$ ,  $x_n$  is obtained by moving on the direction negative gradient of  $L(x)$  at  $x_{n-1}$  with some step size.

If the step size is too large, the iterations will diverge.

If the step size is low enough, the iterations will converge to a stationary point that satisfies  $g(x) = 0$  (it can be a local minimum or saddle point).



# Gradient Descent (GD)

## GD Algorithm

1. Start with arbitrary initial point  $\mathbf{x}_0$ ,
2. For  $n = 1, 2, \dots$ , perform:

$$\mathbf{x}_n = \mathbf{x}_{n-1} - \gamma_n \mathbf{g}(\mathbf{x}_{n-1}), \text{ where } \mathbf{g}(\mathbf{x}) = \nabla_{\mathbf{x}} L(\mathbf{x})$$

3. Upon convergence/termination return  $\mathbf{x}_n$

**Termination:** e.g., when  $\|\mathbf{x}_n - \mathbf{x}_{n-1}\|_2 < \epsilon$  or when  $n > \tau$

**Step size**  $\gamma_n$  must be chosen appropriately; too small: slow convergence; too large: divergence.

# GD for LS

For LS,  $g(x) = \nabla_x L(x) = (A^T Ax - A^T y)2$

GD for LS

1. Start with arbitrary initial point  $x_0$ ,
2. For  $n = 1, 2, \dots$ , perform:

$$x_n = x_{n-1} - \gamma_n (A^T Ax_{n-1} - A^T y)2$$

3. Upon convergence/termination return  $x_n$

For LS,  $g(x) = \nabla_x L(x) = (A^T Ax - A^T y)2$

GD for LS

1. Start with arbitrary initial point  $x_0$ ,
2. For  $n = 1, 2, \dots$ , perform:

$$x_n = x_{n-1} - \gamma_n (A^T Ax_{n-1} - A^T y)2$$

3. Upon convergence/termination return  $x_n$



## GD Step Size – Fixed Step Size

For all  $n$ ,

$$\mathbf{x}_n = \mathbf{x}_{n-1} - \gamma g(\mathbf{x}_{n-1})$$

That is, we keep constant step size for all updates.

Step size determines convergence or divergence.

If  $\gamma$  is too large, the iterations will diverge.

For convergence, constant step size must be **low enough at each step!**

That is, to guarantee convergence, step size must be upper bounded:

$$\gamma < T$$

where  $T$  captures the maximum rate of change of the gradient across the entire domain.

## GD Step Size – Fixed Step Size (Theory)

Let  $L(x)$  have a **Lipschitz Continuous (LC) gradient** with **Lipschitz constant**  $C$  being the minimum positive number for which:

$$\|g(x) - g(y)\| \leq C \|x - y\| \quad \forall x, y$$

That is, the magnitude of the gradient change between any two points is bounded by the distance of these two points, weighted by  $C$ . For the Lipschitz constant it holds:

$$C = \sup_{x, y \in F} \frac{\|g(x) - g(y)\|}{\|x - y\|}$$

Then, if  $L(x)$  is smooth,  $g(x)$  is LC, and we set  $\gamma < 1/C$ , then GD **converges** to a stationary point that satisfies  $g(x) = 0$ .

If  $L(x)$  is also convex the convergence will be at global optimum.

Clearly, high  $C$  means that gradient changes fast; thus, lower  $\gamma$  is needed for convergence.

## GD Step Size – Fixed Step Size (Approximation)

Since it is practically impossible to compute the supremum over all pairs of points in  $F$ , we find a lower bound for  $C$  by examining only the points in a finite search set  $S \subset F$ .

We set  $\gamma < \frac{1}{\hat{C}}$  where

$$\hat{C} = \sup_{x,y \in S} \frac{\|g(x) - g(y)\|}{\|x - y\|}$$

and  $S \subset F$  is a search set.

This approximation comes with no guarantees of convergence –but can work well in practice for appropriate sample set  $S$ .

This approximation can also work when the function is not Lipschitz continuous.

# GD Step Size – Fixed Step Size (C2)

## Special Case: Twice Differentiable Functions

If  $L(x)$  is C2 with LC gradient, the Lipschitz constant of the gradient  $C$  can be found by SVD of the Hessian matrix  $H(x) = \nabla_{\mathbf{x}^T} \nabla_{\mathbf{x}} f(x) = \nabla_{\mathbf{x}^T} g(x)$ . Specifically,  $C$  is given by

$$C = \sup_{x \in F} \sigma_{\max}(H(x))$$

In general, it is practically intractable to compute the supremum over all points in  $F$ .

Instead, we can find a lower bound for  $C$  by evaluating the largest eigenvalue of the Hessian at points in a finite search set  $S \subset F$ . We set  $T < \frac{1}{\hat{C}}$ , where:

$$\hat{C} = \max_{x \in S} \sigma_{\max}(H(x))$$

## GD Step Size – Fixed Step Size (C2)

Even More Special Case: Quadratic Functions: For quadratic functions, the Hessian  $H(\mathbf{x}) = H$  is constant across all points  $\mathbf{x}$  in the domain. Thus,  $C = \sigma_{\max}(H)$ .

For example, for LS objective,

$$H = \nabla_{\mathbf{x}} g(\mathbf{x}) = \nabla_{\mathbf{x}} (\mathbf{A}^T \mathbf{A} \mathbf{x} - \mathbf{A}^T \mathbf{y})^2 = \mathbf{A}^T \mathbf{A} 2$$

Considering RSVDs  $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$  and  $\mathbf{A}^T \mathbf{A} 2 = \mathbf{V} 2 \mathbf{\Sigma}^2 \mathbf{V}^T$ , we find that

$$C = 2 \sigma_{\max}(\mathbf{A})^2$$

# GD Step Size – Gradient Lipschitz Step Size

If  $L(x)$  is  $C^2$  with LC gradient, we can perform gradient descent

$$\mathbf{x}_n = \mathbf{x}_{n-1} - \gamma_n g(\mathbf{x}_{n-1})$$

with step size  $\gamma_n = 1/C_n$ , where

$$C_n = \sigma_{\max} \left( \mathbf{H}(\mathbf{x}_{(n-1)}) \right)$$

is the Lipschitz constant of  $g(\mathbf{x})$  at  $\mathbf{x}_{(n-1)}$ .

Again, this guarantees convergence to a stationary point.

# GD Step Size – Backtracking Line Search

1. fix a parameter  $0 < \beta < 1$  (e.g.,  $\beta = 0.5$ ) and  $c > 0$  (e.g.,  $c = 10^{-2}$ )
2. initialize with arbitrary  $\gamma_n^{cand}$
3. while  $L(x_{n-1}) - L(x_{n-1} - \gamma_n^{cand} g(x_{n-1})) \leq c \cdot \gamma_n^{cand} \|g(x_{n-1})\|^2$  (Armijo Condition)
4.       update  $\gamma_n^{cand} \leftarrow \beta \cdot \gamma_n^{cand}$
5. return  $\gamma_n = \gamma_n^{cand}$

For a smooth function with LC gradient, backtracking line search using Armijo's condition guarantees convergence to a stationary point.

See convergence proofs here: [https://www.cs.cmu.edu/~ggordon/10725-F12/scribes/10725\\_Lecture5.pdf](https://www.cs.cmu.edu/~ggordon/10725-F12/scribes/10725_Lecture5.pdf)

# GD Step Size – Exact Line Search

Step size can be optimized such that

$$\gamma_n = \operatorname{argmin}_{\gamma \geq 0} L(\mathbf{x}_n)$$

where  $\mathbf{x}_n = \mathbf{x}_{n-1} - \gamma \mathbf{g}(\mathbf{x}_{n-1})$  and  $\mathbf{g}(\mathbf{x}_{n-1}) = (A^T A \mathbf{x}_{n-1} - A^T \mathbf{y})/2$

We define  $\mathbf{e}_n := \mathbf{y} - A\mathbf{x}_n \Rightarrow L(\mathbf{x}_n) = \|\mathbf{e}_n\|^2$

And we find that  $\mathbf{e}_n = \mathbf{y} - A(\mathbf{x}_{n-1} - \gamma \mathbf{g}(\mathbf{x}_{n-1})) = \mathbf{y} - A\mathbf{x}_{n-1} + A\mathbf{g}(\mathbf{x}_{n-1})\gamma = \mathbf{e}_{n-1} + A\mathbf{g}(\mathbf{x}_{n-1})\gamma$

We also find that  $\mathbf{g}(\mathbf{x}_{n-1}) = (A^T A \mathbf{x}_{n-1} - A^T \mathbf{y})/2 = -A^T \mathbf{e}_{n-1}/2$

From the above two,  $\mathbf{e}_n = \mathbf{e}_{n-1} + A(-A^T \mathbf{e}_{n-1}/2)\gamma = \mathbf{e}_{n-1} - AA^T \mathbf{e}_{n-1} \gamma$

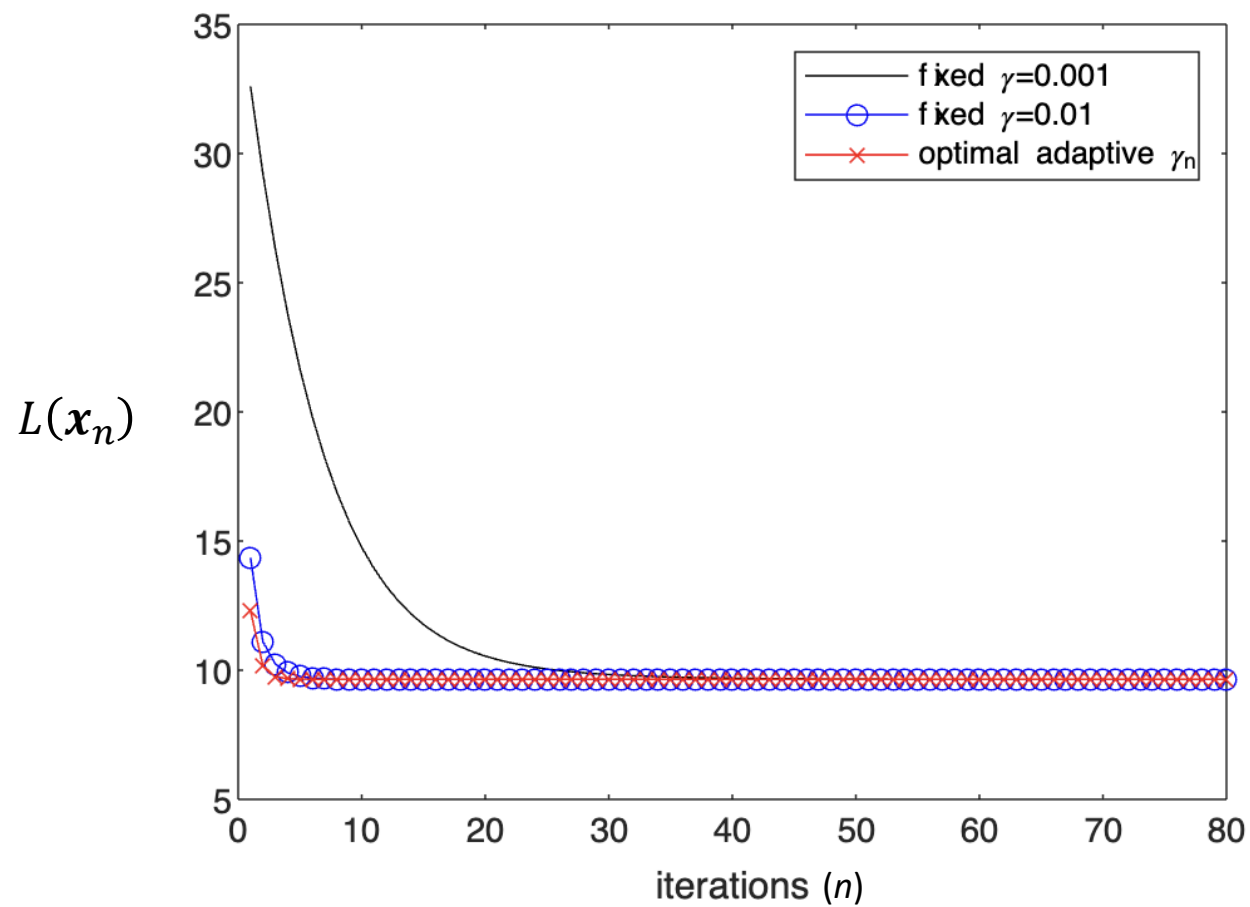
Accordingly,  $L(\mathbf{x}_n) = \|\mathbf{e}_n\|^2 = \|\mathbf{e}_{n-1}\|^2 - 4\gamma \mathbf{e}_{n-1}^T AA^T \mathbf{e}_{n-1} + 4\gamma^2 \|AA^T \mathbf{e}_{n-1}\|^2$

Last, we solve the Exact Line Search as:  $\left. \frac{\partial}{\partial \gamma} L(\mathbf{x}_n) \right|_{\gamma_n} = 0 \Rightarrow \gamma_n = \frac{\|\mathbf{e}_{n-1}\|^2}{2\|AA^T \mathbf{e}_{n-1}\|^2} \Rightarrow \gamma_n = \frac{\|\mathbf{g}(\mathbf{x}_{n-1})\|^2}{2\|A\mathbf{g}(\mathbf{x}_{n-1})\|^2}$



# GD for LS – Step Size (Example)

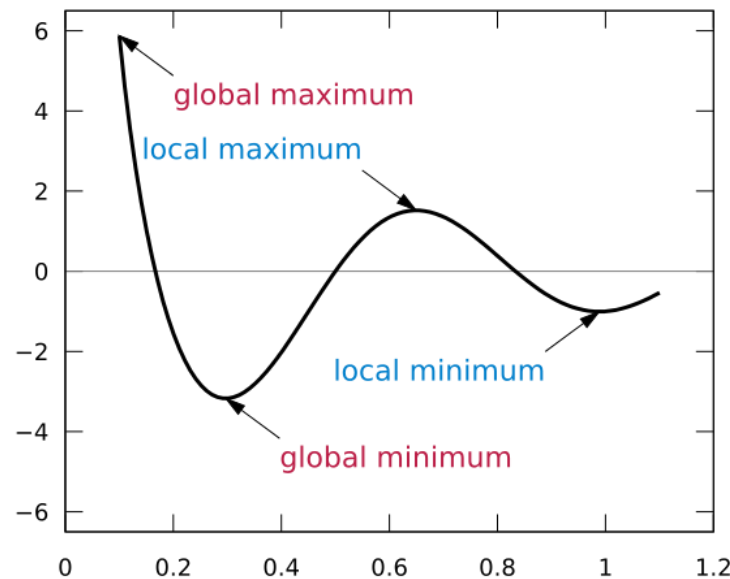
$$A \in \mathbb{R}^{100 \times 10}$$



# GD for LS – Convergence Guarantees

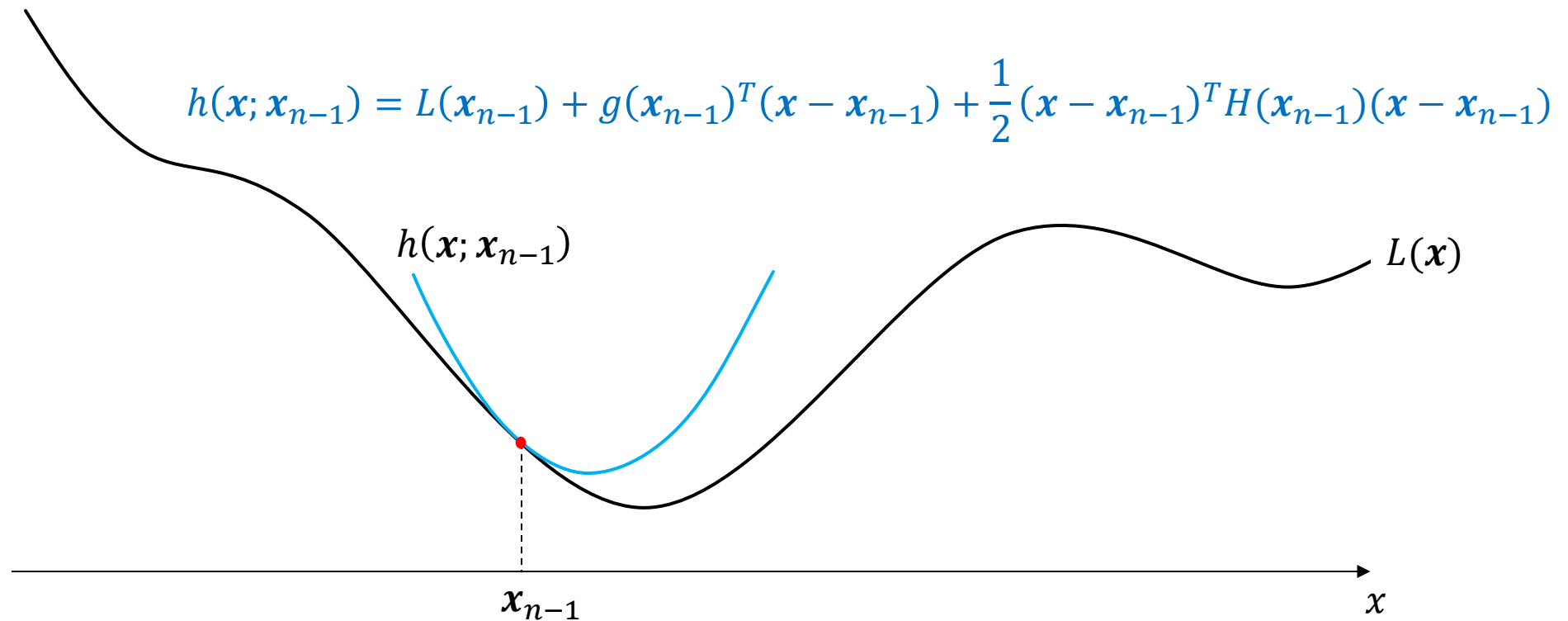
If the function is convex and step sizes appropriately chosen, GD guarantees convergence to **global** minimum  $\mathbf{x}_{opt}$ .

If the function is non-convex, GD converges to a **stationary point**  $\hat{\mathbf{x}}$ . Convergence to the global minimum is possible, subject to “lucky” initial guess, but not guaranteed by GD.



# Newton's Method

- Iterative steps are taken based on the **2<sup>nd</sup> order** Taylor series approximation.
- The minimum is sought by traveling downward on the **quadratics** tangent to  $L(x)$  at each iteration.



# Newton's Method

- Given  $\mathbf{x}_{n-1}$ , set  $\mathbf{x}_n$  to the value of  $\mathbf{x}$  that minimizes

$$h(\mathbf{x}; \mathbf{x}_{n-1}) = L(\mathbf{x}_{n-1}) + g(\mathbf{x}_{n-1})^T(\mathbf{x} - \mathbf{x}_{n-1}) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_{n-1})^T H(\mathbf{x}_{n-1})(\mathbf{x} - \mathbf{x}_{n-1})$$

- For any given  $\mathbf{x}_{n-1}$ ,  $h(\mathbf{x}; \mathbf{x}_{n-1})$  is quadratic over  $\mathbf{x}$  (i.e., convex).
- Minimize over  $\mathbf{x}$  by setting gradient to 0.

$$c(\mathbf{x}; \mathbf{v}) = \nabla h(\mathbf{x}; \mathbf{x}_{n-1}) = g(\mathbf{x}_{n-1}) + H(\mathbf{x}_{n-1})(\mathbf{x} - \mathbf{x}_{n-1})$$

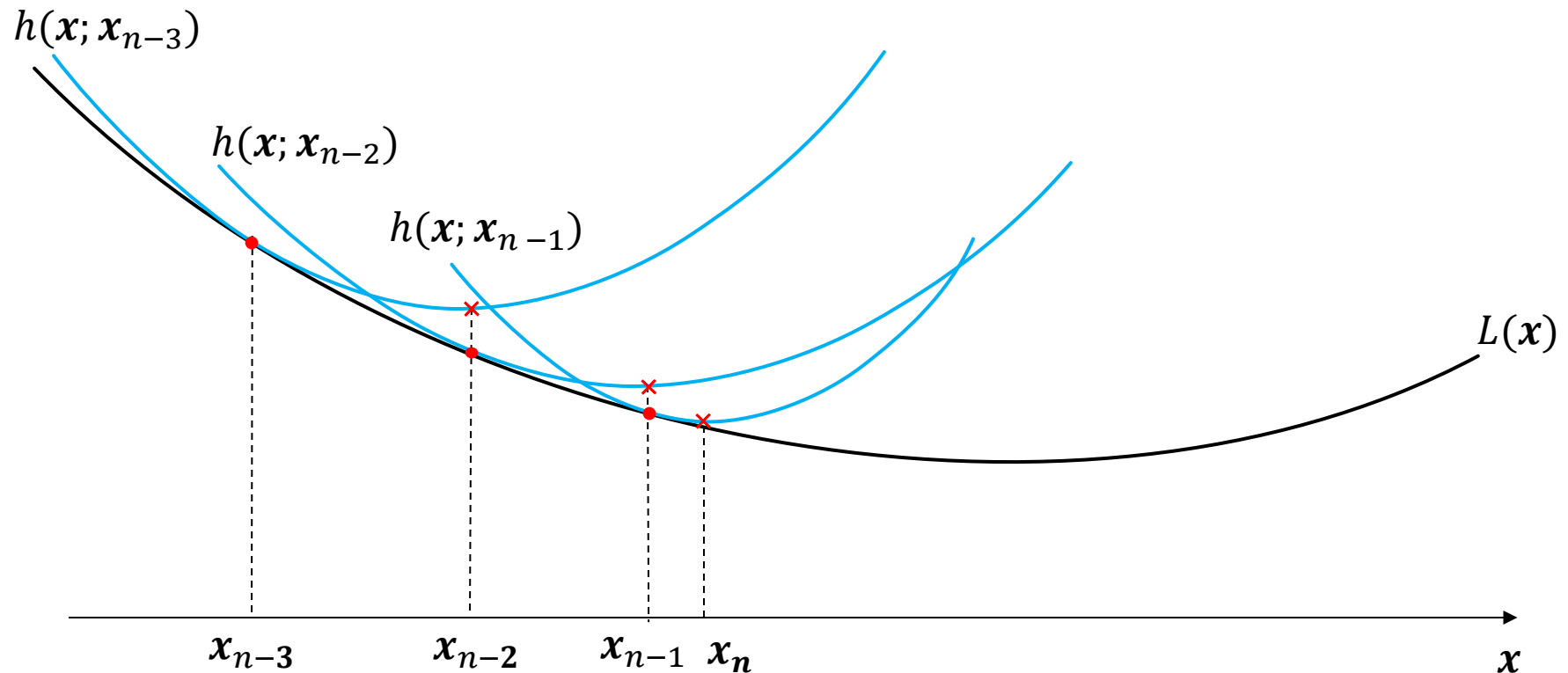
$$c(\mathbf{x}; \mathbf{x}_{n-1})|_{\mathbf{x}=\mathbf{x}_n} = g(\mathbf{x}_{n-1}) + H(\mathbf{v})(\mathbf{x}_n - \mathbf{v}) = 0$$

- Solve the SLE:

$$H(\mathbf{x}_{n-1})\mathbf{x}_n = H(\mathbf{x}_{n-1})\mathbf{x}_{n-1} - g(\mathbf{x}_{n-1})$$

# Newton's Method

1. Begin at arbitrary point  $x_0$
2. For  $n = 1, 2, \dots$  :
  3. Calculate  $H(x_{n-1})$  and  $g(x_{n-1})$
  4. Obtain  $x_n$  that solves the SLE:  $H(x_{k-1})x_k = H(x_{k-1})x_{k-1} - g(x_{k-1})$
5. Upon convergence/termination, return  $x_n$



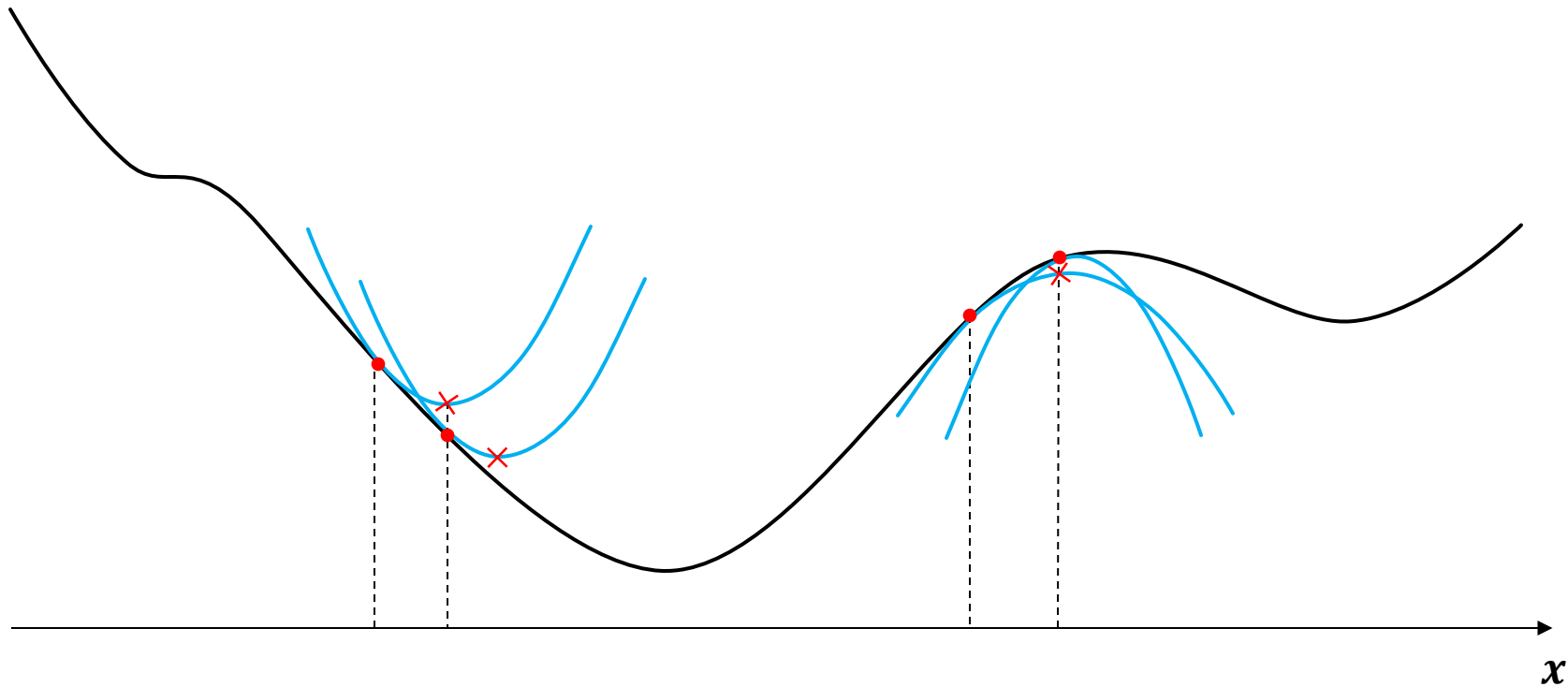
# Newton's Method

For a convex problem, this will converge to global optimum.

For quadratic problem, solutions is obtained at first iteration.

For a non-convex function quadratics can be **concave**.

With such functions Newton's method can *climb* to a maximum.



## LS uses a simple loss...

So far, we have studied the LS problem where the loss  $L(x)$  is a quadratic.

It has a standard closed form solution via SVD or (pseudo)-inversion.

Also, it is a convex problem, and GD returns the global optimum, upon appropriate step size.

Moreover, it has a single Hessian and Lipschitz constant across all values of  $x$ .

Since it is a quadratic, Newton's method solves the problem in a single iteration.

Next, let us look into a different loss function.

It is called “cross-entropy” loss and finds wide application in machine learning. We will see more about its connection to machine learning later in the course. For now, we study it a general optimization problem.

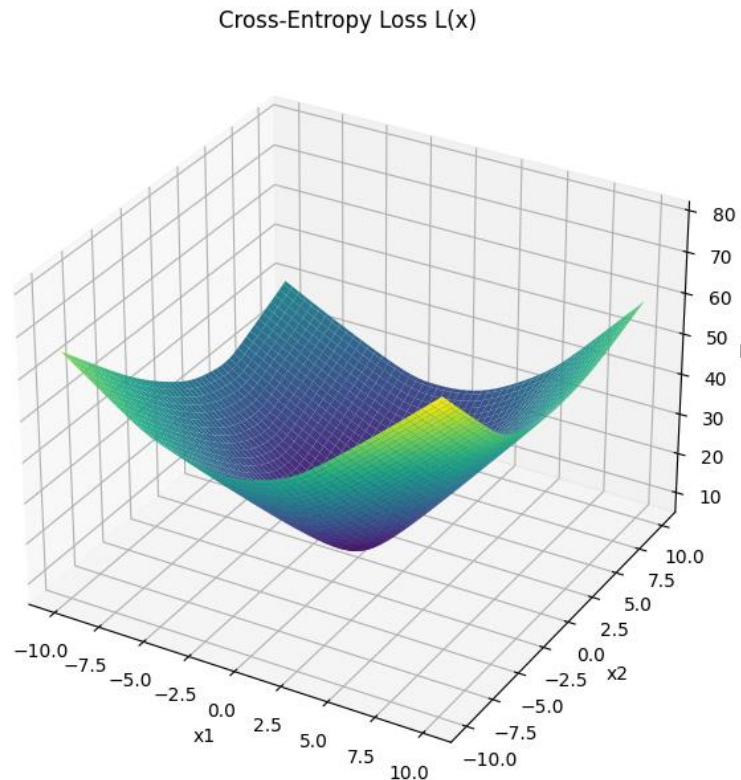
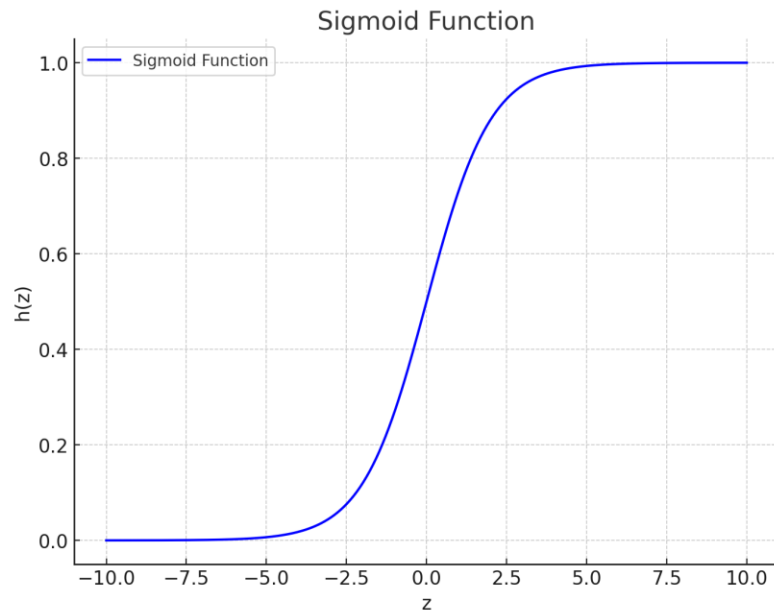
# Cross-Entropy Minimization - Definition

Given  $A \in R^{N \times D}$  and  $\mathbf{y} \in \{0,1\}^N$ , we define the CE objective:

$$L(\mathbf{x}) = - \sum_{i=1}^N \left[ y_i \ln \left( h(A_{i,:} \mathbf{x}) \right) + (1 - y_i) \ln \left( 1 - h(A_{i,:} \mathbf{x}) \right) \right]$$

where  $A_{i,:}$  is the  $i$ -th row of matrix  $A$  and  $h(z) = \frac{1}{1 + \exp(-z)}$  is the **sigmoid** function.

$L(\mathbf{x})$  is still convex but not a quadratic.





# Cross-Entropy Minimization - Gradient

$$\begin{aligned} g(\mathbf{x}) &:= \nabla_{\mathbf{x}} L(\mathbf{x}) = \nabla_{\mathbf{x}} \left( - \sum_{i=1}^N \left[ y_i \ln(h(A_{i,:} \mathbf{x})) + (1 - y_i) \ln(1 - h(A_{i,:} \mathbf{x})) \right] \right) \\ &= - \sum_{i=1}^N \nabla_{\mathbf{x}} \left[ y_i \ln(h(A_{i,:} \mathbf{x})) + (1 - y_i) \ln(1 - h(A_{i,:} \mathbf{x})) \right] = - \sum_{i=1}^N z_i(\mathbf{x}) \end{aligned}$$

where

$$z_i(\mathbf{x}) := \nabla_{\mathbf{x}} \left[ y_i \ln(h(A_{i,:} \mathbf{x})) + (1 - y_i) \ln(1 - h(A_{i,:} \mathbf{x})) \right]$$

We find

$$\begin{aligned} z_i(\mathbf{x}) &= \nabla_{\mathbf{x}} \left[ y_i \ln(h(A_{i,:} \mathbf{x})) + (1 - y_i) \ln(1 - h(A_{i,:} \mathbf{x})) \right] \\ &= y_i \nabla_{\mathbf{x}} (\ln(h(A_{i,:} \mathbf{x}))) + (1 - y_i) \nabla_{\mathbf{x}} (\ln(1 - h(A_{i,:} \mathbf{x}))) \end{aligned}$$

By Chain Rule, for any function  $m(\mathbf{x})$ ,

$$\nabla_{\mathbf{x}} \ln(m(\mathbf{x})) = \frac{1}{m(\mathbf{x})} \nabla_{\mathbf{x}} m(\mathbf{x})$$

and

$$\nabla_{\mathbf{x}} h(m(\mathbf{x})) = h(m(\mathbf{x})) (1 - h(m(\mathbf{x}))) \nabla_{\mathbf{x}} m(\mathbf{x})$$

# Cross-Entropy Minimization - Gradient

Therefore,

$$\begin{aligned}\nabla_{\mathbf{x}} \left( \ln \left( h(A_{i,:}, \mathbf{x}) \right) \right) &= \frac{1}{h(A_{i,:}, \mathbf{x})} h(A_{i,:}, \mathbf{x}) \left( 1 - h(A_{i,:}, \mathbf{x}) \right) A_{i,:}^T \\ &= \left( 1 - h(A_{i,:}, \mathbf{x}) \right) A_{i,:}^T\end{aligned}$$

$$\begin{aligned}\nabla_{\mathbf{x}} \left( \ln \left( 1 - h(A_{i,:}, \mathbf{x}) \right) \right) &= \frac{1}{1 - h(A_{i,:}, \mathbf{x})} \left( - \left( \nabla_{\mathbf{x}} h(A_{i,:}, \mathbf{x}) \right) \right) \\ &= - \frac{1}{1 - h(A_{i,:}, \mathbf{x})} h(A_{i,:}, \mathbf{x}) \left( 1 - h(A_{i,:}, \mathbf{x}) \right) A_{i,:}^T = -h(A_{i,:}, \mathbf{x}) A_{i,:}^T\end{aligned}$$

Accordingly,

$$\begin{aligned}z_i(\mathbf{x}) &= y_i \left( 1 - h(A_{i,:}, \mathbf{x}) \right) A_{i,:}^T + (1 - y_i) \left( -h(A_{i,:}, \mathbf{x}) A_{i,:}^T \right) \\ &= \left( y_i - y_i h(A_{i,:}, \mathbf{x}) + y_i h(A_{i,:}, \mathbf{x}) - h(A_{i,:}, \mathbf{x}) \right) A_{i,:}^T \\ &= \left( y_i - h(A_{i,:}, \mathbf{x}) \right) A_{i,:}^T\end{aligned}$$

# Cross-Entropy Minimization – Gradient, Hessian, Lipschitz

We found  $z_i(\mathbf{x}) = (y_i - h(A_{i,:} \mathbf{x})) A_{i,:}^T, \forall i$ . Therefore, the **CE gradient** is

$$g(\mathbf{x}) = - \sum_{i=1}^N z_i(\mathbf{x}) = - \sum_{i=1}^N (y_i - h(A_{i,:} \mathbf{x})) A_{i,:}^T = \mathbf{A}^T (h(\mathbf{A}\mathbf{x}) - \mathbf{y})$$

Regarding the Hessian,  $H(\mathbf{x}) = \nabla_{\mathbf{x}^T} g(\mathbf{x}) = \nabla_{\mathbf{x}^T} \mathbf{A}^T (h(\mathbf{A}\mathbf{x}) - \mathbf{y}) = \mathbf{A}^T \nabla_{\mathbf{x}^T} h(\mathbf{A}\mathbf{x}) = \mathbf{A}^T \begin{bmatrix} \nabla_{\mathbf{x}^T} h(A_{1,:} \mathbf{x}) \\ \vdots \\ \nabla_{\mathbf{x}^T} h(A_{N,:} \mathbf{x}) \end{bmatrix} = \mathbf{A}^T \begin{bmatrix} h(A_{1,:} \mathbf{x}) (1 - h(A_{1,:} \mathbf{x})) A_{1,:} \\ \vdots \\ h(A_{N,:} \mathbf{x}) (1 - h(A_{N,:} \mathbf{x})) A_{N,:} \end{bmatrix}$

(notice that  $\nabla_{\mathbf{x}} \mathbf{a}^T \mathbf{x} = \mathbf{a}$  and  $\nabla_{\mathbf{x}^T} \mathbf{a}^T \mathbf{x} = \mathbf{a}^T$ )

Therefore, the **CE Hessian** is

$$H(\mathbf{x}) = \mathbf{A}^T R(\mathbf{x}) \mathbf{A}$$

where

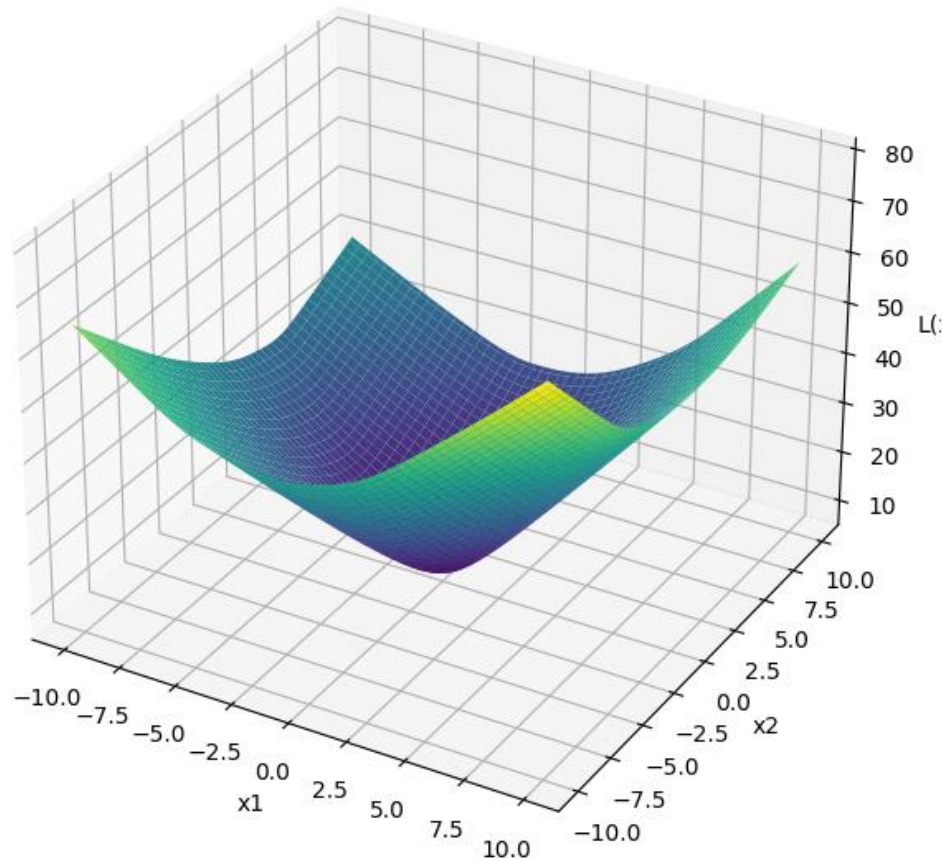
$$R(\mathbf{x}) = \text{diag} \left( \left[ h(A_{1,:} \mathbf{x}) (1 - h(A_{1,:} \mathbf{x})) \right], \dots, \left[ h(A_{N,:} \mathbf{x}) (1 - h(A_{N,:} \mathbf{x})) \right] \right)$$

The Lipschitz Constant of  $L(\mathbf{x})$  at each point  $\mathbf{x}$  is

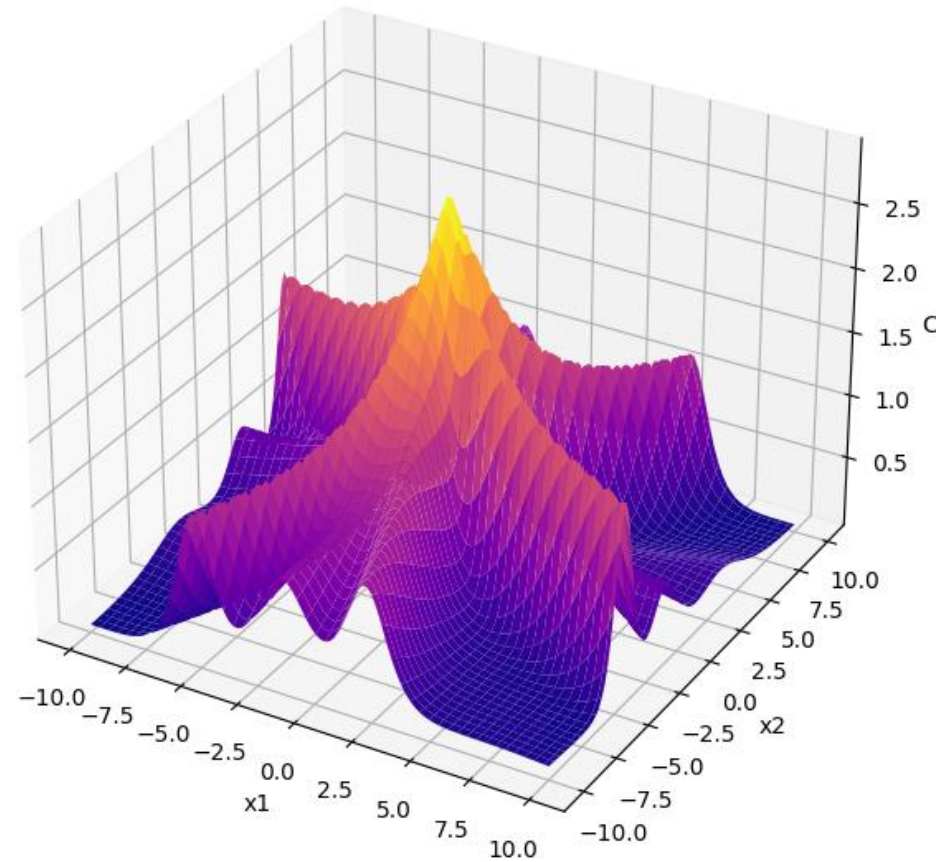
$$C(\mathbf{x}) = \sigma_{\max}(\mathbf{A}^T R(\mathbf{x}) \mathbf{A})$$

# Cross-Entropy Minimization – Lipschitz

Cross-Entropy Loss  $L(x)$



Lipschitz Constant  $C(x)$



# Cross-Entropy Minimization

How could we minimize CE?

No closed form minimizer. No closed form solution to system  $g(\mathbf{x}) = \mathbf{A}^T(h(\mathbf{A}\mathbf{x}) - \mathbf{y}) = \mathbf{0}$

Can solve as:

- 1) GD with FSS. Set step size below  $\frac{1}{C_{min}}$ , where  $C_{min} = \min_{\mathbf{x} \in R^D} C(\mathbf{x}) = \min_{\mathbf{x} \in R^D} \sigma_{max}(H(\mathbf{x}))$ .

Can use instead  $\hat{C}_{min} = \min_{\mathbf{x} \in S} \sigma_{max}(H(\mathbf{x}))$  where  $S$  is a subset of  $R^D$

- 2) GD with ELS
- 3) GD with BLS
- 4) GD with Lipschitz constant steps; i.e.,  $\gamma_n = \sigma_{max}(H(\mathbf{x}_{n-1}))$
- 5) Newton's method (non-trivial, multiple iterations in contrast to LS)

...

and many more approaches such as stochastic variants of GD that we will talk about later on.