

# Basics of Numerical Optimization: Optimality Conditions

---

**Ju Sun**

Computer Science & Engineering  
University of Minnesota, Twin Cities

September 20, 2023

- HW0 is out (deadline: Sep 29th, i.e., 10 days)
  - \* ChatGPT, Claude, Github Copilot, and so on are your friends!
  - \* We'll demo a bit of this in class
- Start playing with Google Colab (<https://colab.google/>)  
[Reminder: CS department pays 3-month Colab Pro subscription into your student account]
- MSI class account is on the way

# Supervised learning as data fitting

Step	General view	NN view
1	Gather training set $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$	Gather training set $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)$
2	Choose a family of functions, e.g., $\mathcal{H}$ , so that there is an $f \in \mathcal{H}$ to ensure $\mathbf{y}_i \approx f(\mathbf{x}_i), \forall i$	Choose a NN with $k$ neurons, so that there is a group of weights $(\mathbf{w}_1, \dots, \mathbf{w}_k, b_1, \dots, b_k)$ ensuring $\mathbf{y}_i \approx \{\text{NN}(\mathbf{w}_1, \dots, \mathbf{w}_k, b_1, \dots, b_k)\}(\mathbf{x}_i), \forall i$
3	Set up a loss function $\ell$	Set up a loss function $\ell$
4	Find an $f \in \mathcal{H}$ to minimize the average loss $\frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}_i, f(\mathbf{x}_i))$	Find weights $(\mathbf{w}_1, \dots, \mathbf{w}_k, b_1, \dots, b_k)$ to minimize the average loss $\frac{1}{n} \sum_{i=1}^n \ell[\mathbf{y}_i, \{\text{NN}(\mathbf{w}_1, \dots, \mathbf{w}_k, b_1, \dots, b_k)\}(\mathbf{x}_i)]$

## Three fundamental questions in DL

- **Approximation:** is it powerful, i.e., the  $\mathcal{H}$  large enough for all possible weights? (last lecture: universal approximation theorems)
  - \* two-layer neural networks can be universal approximators
  - \* deep neural networks more efficient than shallow neural networks
- **Optimization:** how to solve

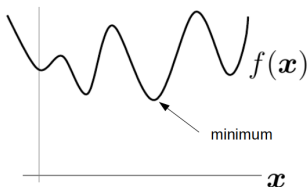
$$\min_{\mathbf{w}'_i s, \mathbf{b}'_i s} \frac{1}{n} \sum_{i=1}^n \ell[\mathbf{y}_i, \{\text{NN}(\mathbf{w}_1, \dots, \mathbf{w}_k, b_1, \dots, b_k)\}(\mathbf{x}_i)]$$

(now)

- **Generalization:** does the learned NN work well on “similar” data? (CSCI5525, and Deep Learning Theory)

Optimality conditions of unconstrained optimization

# Optimization problems



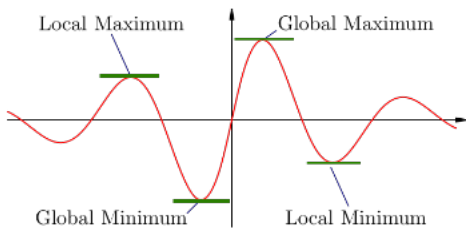
*Nothing takes place in the world whose meaning is not that of some maximum or minimum. – Euler*

$$\min_x f(x) \text{ s. t. } x \in C.$$

- $x$ : optimization variables,  $f(x)$ : objective function,  $C$ : constraint (or feasible) set
- $C$  consists of continuous values (e.g.,  $\mathbb{R}^n$ ,  $[0, 1]^n$ ): **continuous optimization**;  $C$  consists of discrete values (e.g.,  $\{-1, +1\}^n$ ): discrete optimization
- $C$  whole space  $\mathbb{R}^n$ : **unconstrained optimization**;  $C$  a strict subset of the space: constrained optimization

We focus on **continuous, unconstrained** optimization here.

# Global and local mins



Credit: study.com

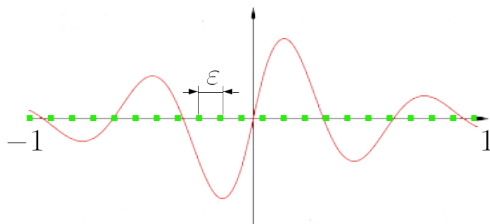
Let  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ ,

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

- $\mathbf{x}_0$  is a **local minimizer** if:  $\exists \varepsilon > 0$ , so that  $f(\mathbf{x}_0) \leq f(\mathbf{x})$  for all  $\mathbf{x}$  satisfying  $\|\mathbf{x} - \mathbf{x}_0\|_2 < \varepsilon$ . The value  $f(\mathbf{x}_0)$  is called a **local minimum**.
- $\mathbf{x}_0$  is a **global minimizer** if:  $f(\mathbf{x}_0) \leq f(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}^n$ . The value is  $f(\mathbf{x}_0)$  called **the global minimum**.

# A naive method for optimization

## Grid search

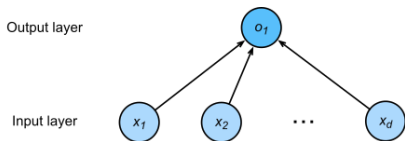


- For 1D problem, assume we know the global min lies in  $[-1, 1]$
- Take uniformly grid points in  $[-1, 1]$  so that any adjacent points are separated by  $\epsilon$ .
- Need  $O(\epsilon^{-1})$  points to get an  $\epsilon$ -close point to the global min by exhaustive search

For  $N$ -D problems, need  $O(\epsilon^{-n})$  computation.



# What we do in practice



Credit: D2L

$\sigma$  is the identity function

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \|y_i - \mathbf{w}^\top \mathbf{x}_i\|_2^2$$

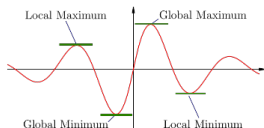
$$\min_{\mathbf{w}} f(\mathbf{w}) \doteq \frac{1}{n} \sum_{i=1}^n \|y_i - \mathbf{w}^\top \mathbf{x}_i\|_2^2 = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad \text{where } \mathbf{X} \doteq \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix}$$
$$\implies \nabla f(\mathbf{w}) = \frac{2}{n} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$\nabla f(\mathbf{w}) = \mathbf{0} \iff \frac{2}{n} \mathbf{X}^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{0} \implies \mathbf{w} = (\mathbf{X}^\top \mathbf{X})^+ \mathbf{X}^\top \mathbf{y} + \text{null}(\mathbf{X})$$

**Optimality conditions:** Reduce the search space by **characterizing** the local/global mins

# First-order optimality condition

**Necessary condition:** Assume  $f$  is 1st-order differentiable at  $\mathbf{x}_0$ . If  $\mathbf{x}_0$  is a local minimizer,  $\nabla f(\mathbf{x}_0) = \mathbf{0}$ .



**Intuition:**  $\nabla f$  is “rate of change” of function value. If the rate is not zero at  $\mathbf{x}_0$ , possible to decrease  $f$  along  $-\nabla f(\mathbf{x}_0)$

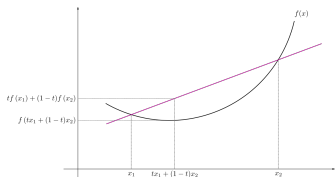
Taylor's:  $f(\mathbf{x}_0 + \boldsymbol{\delta}) = f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2)$ . If  $\mathbf{x}_0$  is a local min:

- For all  $\boldsymbol{\delta}$  sufficiently small,  
 $f(\mathbf{x}_0 + \boldsymbol{\delta}) - f(\mathbf{x}_0) = \langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2) \geq 0$
- For all  $\boldsymbol{\delta}$  sufficiently small, sign of  $\langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2)$  determined by the sign of  $\langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle$ , i.e.,  $\langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle \geq 0$ .
- So for all  $\boldsymbol{\delta}$  sufficiently small,  $\langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle \geq 0$  and  
 $\langle \nabla f(\mathbf{x}_0), -\boldsymbol{\delta} \rangle = -\langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle \geq 0 \implies \langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle = 0$
- So  $\nabla f(\mathbf{x}_0) = \mathbf{0}$ .

# First-order optimality condition

**Necessary condition:** Assume  $f$  is 1st-order differentiable at  $\mathbf{x}_0$ . If  $\mathbf{x}_0$  is a local minimizer, then  $\nabla f(\mathbf{x}_0) = \mathbf{0}$ .

When sufficient? **for convex functions**



Credit: Wikipedia

- **geometric def.:** function for which any line segment connecting two points of its graph always lies above the graph
- **algebraic def.:**  $\forall \mathbf{x}, \mathbf{y}$  and  $\alpha \in [0, 1]$ :

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}).$$

Any convex function has only **one** local minimum (**value!**), **which is also global!**

**Proof sketch:** if  $\mathbf{x}, \mathbf{z}$  are both local minimizers and  $f(\mathbf{z}) < f(\mathbf{x})$ ,

$$f(\alpha \mathbf{z} + (1 - \alpha) \mathbf{x}) \leq \alpha f(\mathbf{z}) + (1 - \alpha) f(\mathbf{x}) < \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{x}) = f(\mathbf{x}).$$

But  $\alpha \mathbf{z} + (1 - \alpha) \mathbf{x} \rightarrow \mathbf{x}$  as  $\alpha \rightarrow 0$ .

# First-order optimality condition

**Necessary condition:** Assume  $f$  is 1st-order differentiable at  $\mathbf{x}_0$ . If  $\mathbf{x}_0$  is a local minimizer, then  $\nabla f(\mathbf{x}_0) = \mathbf{0}$ .

**Sufficient condition:** Assume  $f$  is **convex and 1st-order differentiable**. If  $\nabla f(\mathbf{x}) = \mathbf{0}$  at a point  $\mathbf{x} = \mathbf{x}_0$ , then  $\mathbf{x}_0$  is a local/global minimizer.

- Suppose  $f$  is twice differentiable.  $f$  is convex  $\iff \nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$  for all  $\mathbf{x}$ 
  - \* Consider  $f(\mathbf{w}) = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$  and its solutions again
  - \* Is  $f(\mathbf{W}_1, \mathbf{W}_2) = \|\mathbf{y} - \mathbf{W}_2\mathbf{W}_1\mathbf{x}\|_2^2$  convex?
- Convex analysis (i.e., theory) and optimization (i.e., numerical methods) are relatively mature. Recommended resources: analysis: [Hiriart-Urruty and Lemaréchal, 2001], optimization: [Boyd and Vandenberghe, 2004]
- We **don't assume** convexity unless stated, as DNN objectives are almost always nonconvex.

## Second-order optimality condition

**Necessary condition:** Assume  $f(\mathbf{x})$  is 2-order differentiable at  $\mathbf{x}_0$ . If  $\mathbf{x}_0$  is a local min,  $\nabla f(\mathbf{x}_0) = \mathbf{0}$  and  $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$  (i.e., positive semidefinite).

**Sufficient condition:** Assume  $f(\mathbf{x})$  is 2-order differentiable at  $\mathbf{x}_0$ . If  $\nabla f(\mathbf{x}_0) = \mathbf{0}$  and  $\nabla^2 f(\mathbf{x}_0) \succ \mathbf{0}$  (i.e., positive definite),  $\mathbf{x}_0$  is a local min.

Taylor's:  $f(\mathbf{x}_0 + \boldsymbol{\delta}) = f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}_0) \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2^2)$ .

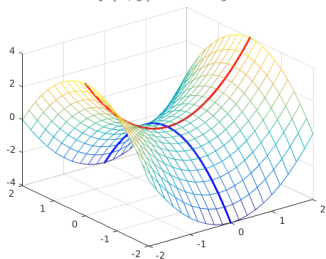
- If  $\mathbf{x}_0$  is a local min,  $\nabla f(\mathbf{x}_0) = \mathbf{0}$  (1st-order condition) and  $f(\mathbf{x}_0 + \boldsymbol{\delta}) = f(\mathbf{x}_0) + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}_0) \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2^2)$ .
- So  $f(\mathbf{x}_0 + \boldsymbol{\delta}) - f(\mathbf{x}_0) = \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}_0) \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2^2) \geq 0$  for all  $\boldsymbol{\delta}$  sufficiently small
- For all  $\boldsymbol{\delta}$  sufficiently small, sign of  $\frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}_0) \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2^2)$  determined by the sign of  $\frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}_0) \boldsymbol{\delta} \rangle \implies \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}_0) \boldsymbol{\delta} \rangle \geq 0$
- So  $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$ .

# What's in between?

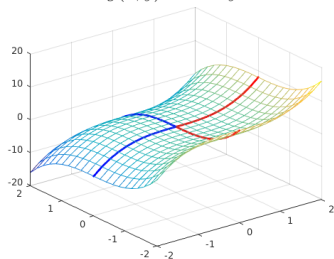
2nd order sufficient:  $\nabla f(\mathbf{x}_0) = \mathbf{0}$  and  $\nabla^2 f(\mathbf{x}_0) \succ \mathbf{0}$

2nd order necessary:  $\nabla f(\mathbf{x}_0) = \mathbf{0}$  and  $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$

$$f(x, y) = x^2 - y^2$$



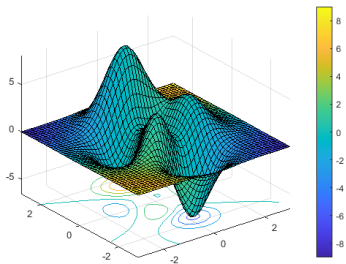
$$g(x, y) = x^3 - y^3$$



$$\nabla f = \begin{bmatrix} 2x \\ -2y \end{bmatrix}, \nabla^2 f = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

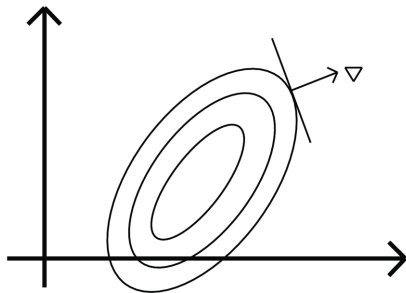
$$\nabla g = \begin{bmatrix} 3x^2 \\ -3y^2 \end{bmatrix}, \nabla^2 g = \begin{bmatrix} 6x & 0 \\ 0 & -6y \end{bmatrix}$$

# Contour plot



contour/levelset plot

(Credit: Mathworks)



gradient direction? why?

- [Boyd and Vandenberghe, 2004] Boyd, S. and Vandenberghe, L. (2004). **Convex Optimization**. Cambridge University Press.
- [Hiriart-Urruty and Lemaréchal, 2001] Hiriart-Urruty, J.-B. and Lemaréchal, C. (2001). **Fundamentals of Convex Analysis**. Springer Berlin Heidelberg.