

# Basics of Numerical Optimization: Preliminaries

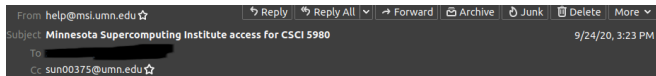
---

**Ju Sun**

Computer Science & Engineering  
University of Minnesota, Twin Cities

September 28, 2020

## – MSI account



This email is to inform you that you have received access to the Minnesota Supercomputing Institute for use in conjunction with your Fall 2020 CSCI 5980 coursework.

Before you can begin using MSI systems, you will need to review our User Agreement by visiting the URL below and logging in with your University of Minnesota Internet ID:

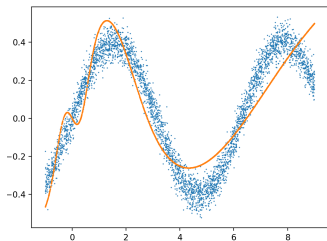
<https://www.msi.umn.edu/user-agreement>

The website will provide MSI login information to you after you agree to the statement.

For assistance in using your account after you have completed this process, you should work with your instructor.

- Oct 05 Zoom-only lecture: Intro to Numpy/Scipy/Colab/MSI
- + Discussion of project ideas

# Supervised learning as function approximation



- Underlying true function:  $f_0$
- Training data:  $\{\mathbf{x}_i, \mathbf{y}_i\}$  with  $\mathbf{y}_i \approx f_0(\mathbf{x}_i)$
- Choose a family of functions  $\mathcal{H}$ , so that  $\exists f \in \mathcal{H}$  and  $f$  and  $f_0$  are close
- Find  $f$ , i.e., optimization

$$\min_{f \in \mathcal{H}} \sum_i \ell(\mathbf{y}_i, f(\mathbf{x}_i)) + \Omega(f)$$

- **Approximation capacity: Universal approximation theorems (UAT)**  
 $\implies$  replace  $\mathcal{H}$  by  $\text{DNN}_{\mathbf{W}}$ , i.e., a deep neural network with weights  $\mathbf{W}$
- **Optimization:**

$$\min_{\mathbf{W}} \sum_i \ell(\mathbf{y}_i, \text{DNN}_{\mathbf{W}}(\mathbf{x}_i)) + \Omega(\mathbf{W})$$

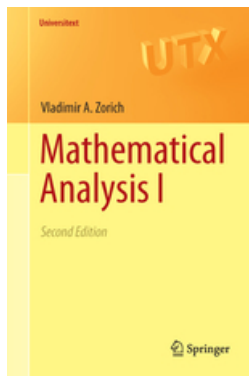
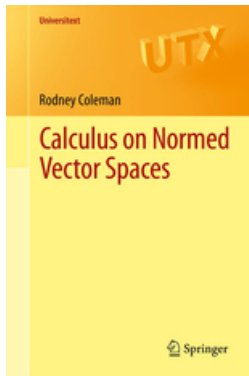
- **Generalization:** how to avoid over-complicated  $\text{DNN}_{\mathbf{W}}$  in view of UAT

Now we start to focus on **optimization**.

Review of high-dimensional calculus

Optimality conditions of unconstrained optimization

## Recommended references



[Munkres, 1997, Coleman, 2012, Zorich, 2015]

## Our notation

- scalars:  $x$ , vectors:  $\mathbf{x}$ , matrices:  $\mathbf{X}$ , tensors:  $\mathcal{X}$ , sets:  $S$
- vectors are always **column vectors**, unless stated otherwise
- $x_i$ :  $i$ -th element of  $\mathbf{x}$ ,  $x_{ij}$ :  $(i, j)$ -th element of  $\mathbf{X}$ ,  $\mathbf{x}^i$ :  $i$ -th row of  $\mathbf{X}$  as a **row vector**,  $\mathbf{x}_j$ :  $j$ -th column of  $\mathbf{X}$  as a **column vector**
- $\mathbb{R}$ : real numbers,  $\mathbb{R}_+$ : positive reals,  $\mathbb{R}^n$ : space of  $n$ -dimensional vectors,  $\mathbb{R}^{m \times n}$ : space of  $m \times n$  matrices,  $\mathbb{R}^{m \times n \times k}$ : space of  $m \times n \times k$  tensors, etc
- $[n] \doteq \{1, \dots, n\}$

# Differentiability — first order

Consider  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Definition: **First-order differentiable** at a point  $\mathbf{x}$  if there exists a matrix  $\mathbf{B} \in \mathbb{R}^{m \times n}$  such that

$$\frac{f(\mathbf{x} + \boldsymbol{\delta}) - f(\mathbf{x}) - \mathbf{B}\boldsymbol{\delta}}{\|\boldsymbol{\delta}\|_2} \rightarrow \mathbf{0} \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

$$\text{i.e., } f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \mathbf{B}\boldsymbol{\delta} + o(\|\boldsymbol{\delta}\|_2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

- $\mathbf{B}$  is called the (Fréchet) derivative. When  $m = 1$ ,  $\mathbf{b}^\top$  (i.e.,  $\mathbf{B}^\top$ ) called **gradient**, denoted as  $\nabla f(\mathbf{x})$ . For general  $m$ , also called **Jacobian** matrix, denoted as  $\mathbf{J}_f(\mathbf{x})$ .
- Calculation:  $b_{ij} = \frac{\partial f_i}{\partial x_j}(\mathbf{x})$
- **Sufficient condition**: if all partial derivatives exist and are **continuous** at  $\mathbf{x}$ , then  $f(\mathbf{x})$  is differentiable at  $\mathbf{x}$ .

# Calculus rules

Assume  $f, g : \mathbb{R}^n \rightarrow \mathbb{R}^m$  are differentiable at a point  $\mathbf{x} \in \mathbb{R}^n$ .

- **linearity:**  $\lambda_1 f + \lambda_2 g$  is differentiable at  $\mathbf{x}$  and
$$\nabla [\lambda_1 f + \lambda_2 g] (\mathbf{x}) = \lambda_1 \nabla f (\mathbf{x}) + \lambda_2 \nabla g (\mathbf{x})$$
- **product:** assume  $m = 1$ ,  $fg$  is differentiable at  $\mathbf{x}$  and
$$\nabla [fg] (\mathbf{x}) = f (\mathbf{x}) \nabla g (\mathbf{x}) + g (\mathbf{x}) \nabla f (\mathbf{x})$$
- **quotient:** assume  $m = 1$  and  $g (\mathbf{x}) \neq 0$ ,  $\frac{f}{g}$  is differentiable at  $\mathbf{x}$  and
$$\nabla \left[ \frac{f}{g} \right] (\mathbf{x}) = \frac{g (\mathbf{x}) \nabla f (\mathbf{x}) - f (\mathbf{x}) \nabla g (\mathbf{x})}{g^2 (\mathbf{x})}$$
- **Chain rule:** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and  $h : \mathbb{R}^n \rightarrow \mathbb{R}^k$ , and  $f$  is differentiable at  $\mathbf{x}$  and  $\mathbf{y} = f (\mathbf{x})$  and  $h$  is differentiable at  $\mathbf{y}$ . Then,  $h \circ f : \mathbb{R}^n \rightarrow \mathbb{R}^k$  is differentiable at  $\mathbf{x}$ , and

$$\mathbf{J}_{[h \circ f]} (\mathbf{x}) = \mathbf{J}_h (f (\mathbf{x})) \mathbf{J}_f (\mathbf{x}).$$

When  $k = 1$ ,

$$\nabla [h \circ f] (\mathbf{x}) = \mathbf{J}_f^\top (\mathbf{x}) \nabla h (f (\mathbf{x})).$$



## Put the definition in good use!

**First-order differentiable** at a point  $x$  if there exists a matrix  $B \in \mathbb{R}^{m \times n}$ , called Jacobian, such that

$$f(x + \delta) = f(x) + B\delta + o(\|\delta\|_2) \quad \text{as } \delta \rightarrow \mathbf{0}.$$

- prove the chain rule for  $h \circ f(x)$  (whiteboard)
- derive Jacobian (white board)

$$* f(x) = Ax$$

$$* g(W_1, W_2, W_3) = y - W_1 W_2 W_3 x$$

## Differentiability — second order

Consider  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  and assume  $f$  is 1st-order differentiable in a small ball around  $\mathbf{x}$

- Write  $\frac{\partial f^2}{\partial x_j \partial x_i}(\mathbf{x}) \doteq \left[ \frac{\partial}{\partial x_j} \left( \frac{\partial f}{\partial x_i} \right) \right](\mathbf{x})$  provided the right side well defined
- **Symmetry:** If both  $\frac{\partial f^2}{\partial x_j \partial x_i}(\mathbf{x})$  and  $\frac{\partial f^2}{\partial x_i \partial x_j}(\mathbf{x})$  exist and both are continuous at  $\mathbf{x}$ , then **they are equal**.
- **Hessian (matrix):**

$$\nabla^2 f(\mathbf{x}) \doteq \left[ \frac{\partial f^2}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i}, \quad (1)$$

where  $\left[ \frac{\partial f^2}{\partial x_j \partial x_i}(\mathbf{x}) \right]_{j,i} \in \mathbb{R}^{n \times n}$  has its  $(j, i)$ -th element as  $\frac{\partial f^2}{\partial x_j \partial x_i}(\mathbf{x})$ .

- $\nabla^2 f$  is symmetric.
- **Sufficient condition:** if all  $\frac{\partial f^2}{\partial x_j \partial x_i}(\mathbf{x})$  exist and are **continuous**,  $f$  is 2nd-order differentiable at  $\mathbf{x}$  (**not converse; we omit the definition due to its technicality**).

# Taylor's theorem

**Vector version:** consider  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$

- If  $f$  is 1st-order differentiable at  $\mathbf{x}$ , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

- If  $f$  is 2nd-order differentiable at  $\mathbf{x}$ , then

$$f(\mathbf{x} + \boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}) \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2^2) \text{ as } \boldsymbol{\delta} \rightarrow \mathbf{0}.$$

**Matrix version:** consider  $f(\mathbf{X}) : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

- If  $f$  is 1st-order differentiable at  $\mathbf{X}$ , then

$$f(\mathbf{X} + \boldsymbol{\Delta}) = f(\mathbf{X}) + \langle \nabla f(\mathbf{X}), \boldsymbol{\Delta} \rangle + o(\|\boldsymbol{\Delta}\|_F) \text{ as } \boldsymbol{\Delta} \rightarrow \mathbf{0}.$$

- If  $f$  is 2nd-order differentiable at  $\mathbf{X}$ , then

$$f(\mathbf{X} + \boldsymbol{\Delta}) = f(\mathbf{X}) + \langle \nabla f(\mathbf{X}), \boldsymbol{\Delta} \rangle + \frac{1}{2} \langle \boldsymbol{\Delta}, \nabla^2 f(\mathbf{X}) \boldsymbol{\Delta} \rangle + o(\|\boldsymbol{\Delta}\|_F^2) \\ \text{as } \boldsymbol{\Delta} \rightarrow \mathbf{0}.$$

## Put Taylor in good use!

- derive gradient and Hessian for  $f(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$   
(whiteboard)

- derive gradient (and Hessian) for

$$g(\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3) = \|\mathbf{y} - \mathbf{W}_1\mathbf{W}_2\mathbf{W}_3\mathbf{x}\|_F^2$$

(whiteboard)

before: **gradient, Hessian**  $\implies$  **Taylor expansion**

now: **Taylor expansion**  $\implies$  **gradient, Hessian**

But why?

## Taylor approximation — asymptotic uniqueness

Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be  $k$  ( $k \geq 1$  integer) times differentiable at a point  $x$ . If  $P(\delta)$  is a  $k$ -th order polynomial satisfying  $f(x + \delta) - P(\delta) = o(\delta^k)$  as  $\delta \rightarrow 0$ , then  $P(\delta) = P_k(\delta) \doteq f(x) + \sum_{i=1}^k \frac{1}{i!} f^{(i)}(x) \delta^i$ .

### Generalization to the vector version

- Assume  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  is 1-order differentiable at  $\mathbf{x}$ . If  $P(\boldsymbol{\delta}) \doteq f(\mathbf{x}) + \langle \mathbf{v}, \boldsymbol{\delta} \rangle$  satisfies that

$$f(\mathbf{x} + \boldsymbol{\delta}) - P(\boldsymbol{\delta}) = o(\|\boldsymbol{\delta}\|_2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0},$$

then  $P(\boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle$ , i.e., the 1st-order Taylor expansion.

- Assume  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$  is 2-order differentiable at  $\mathbf{x}$ . If  $P(\boldsymbol{\delta}) \doteq f(\mathbf{x}) + \langle \mathbf{v}, \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \mathbf{H} \boldsymbol{\delta} \rangle$  with  $\mathbf{H}$  symmetric satisfies that

$$f(\mathbf{x} + \boldsymbol{\delta}) - P(\boldsymbol{\delta}) = o(\|\boldsymbol{\delta}\|_2^2) \quad \text{as } \boldsymbol{\delta} \rightarrow \mathbf{0},$$

then  $P(\boldsymbol{\delta}) = f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \boldsymbol{\delta} \rangle + \frac{1}{2} \langle \boldsymbol{\delta}, \nabla^2 f(\mathbf{x}) \boldsymbol{\delta} \rangle$ , i.e., the 2nd-order Taylor expansion. **We can read off  $\nabla f$  and  $\nabla^2 f$  if we know the expansion!**

**Similarly for the matrix version.** See Chap 5 of [Coleman, 2012] for other forms of Taylor theorems and proofs of the asymptotic uniqueness.

# Directional derivatives and curvatures

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

- **directional derivative:**  $D_{\mathbf{v}}f(\mathbf{x}) \doteq \frac{d}{dt}f(\mathbf{x} + t\mathbf{v})$
- When  $f$  is 1-st order differentiable at  $\mathbf{x}$ ,

$$D_{\mathbf{v}}f(\mathbf{x}) = \langle \nabla f(\mathbf{x}), \mathbf{v} \rangle.$$

- Now  $D_{\mathbf{v}}f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ , what is  $D_{\mathbf{u}}(D_{\mathbf{v}}f)(\mathbf{x})$ ?

$$D_{\mathbf{u}}(D_{\mathbf{v}}f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{v} \rangle.$$

- When  $\mathbf{u} = \mathbf{v}$ ,

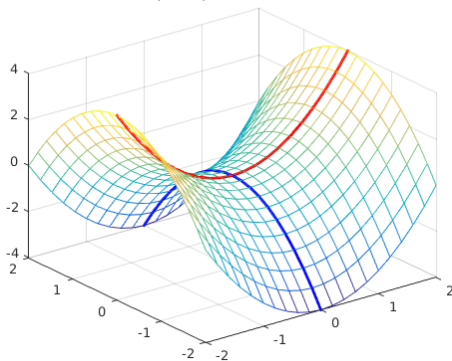
$$D_{\mathbf{u}}(D_{\mathbf{u}}f)(\mathbf{x}) = \langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle = \frac{d^2}{dt^2}f(\mathbf{x} + t\mathbf{u}).$$

- $\frac{\langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle}{\|\mathbf{u}\|_2^2}$  is the **directional curvature** along  $\mathbf{u}$  independent of the norm of  $\mathbf{u}$

# Directional curvature

$\frac{\langle \mathbf{u}, \nabla^2 f(\mathbf{x}) \mathbf{u} \rangle}{\|\mathbf{u}\|_2^2}$  is the **directional curvature** along  $\mathbf{u}$  independent of the norm of  $\mathbf{u}$

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



Blue: negative curvature (bending down)

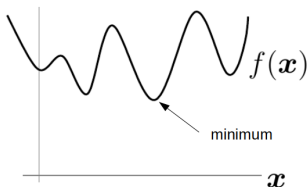
Red: positive curvature (bending up)

Review of high-dimensional calculus

Optimality conditions of unconstrained optimization



# Optimization problems



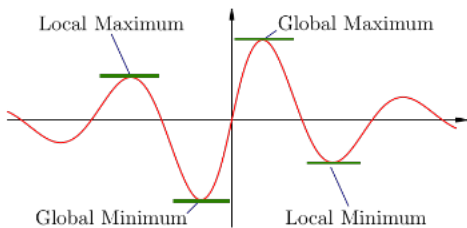
*Nothing takes place in the world whose meaning is not that of some maximum or minimum. – Euler*

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

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

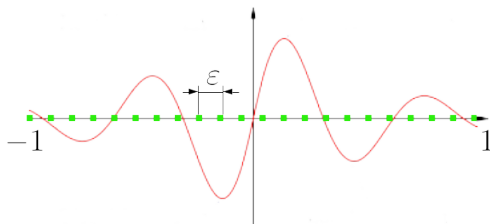
Let  $f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ ,

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

- $\mathbf{x}_0$  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 solution

## Grid search



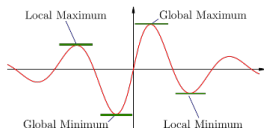
- For 1D problem, assume we know the global min lies in  $[-1, 1]$
- Take uniformly grid points in  $[-1, 1]$  so that any adjacent points are separated by  $\epsilon$ .
- Need  $O(\epsilon^{-1})$  points to get an  $\epsilon$ -close point to the global min by exhaustive search

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

Better characterization of the local/global mins may help avoid this.

# 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  $-\nabla f(\mathbf{x}_0)$

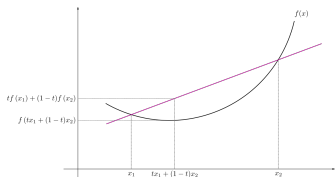
Taylor's:  $f(\mathbf{x}_0 + \boldsymbol{\delta}) = f(\mathbf{x}_0) + \langle \nabla f(\mathbf{x}_0), \boldsymbol{\delta} \rangle + o(\|\boldsymbol{\delta}\|_2)$ . If  $\mathbf{x}_0$  is a local min:

- 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  $x_0$ . If  $x_0$  is a local minimizer, then  $\nabla f(x_0) = \mathbf{0}$ .

**Sufficient condition:** Assume  $f$  is **convex and 1st-order differentiable**. If  $\nabla f(x) = \mathbf{0}$  at a point  $x = x_0$ , then  $x_0$  is a local/global minimizer.

- 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  $f(x)$  is 2-order differentiable at  $x_0$ . If  $x_0$  is a local min,  $\nabla f(x_0) = \mathbf{0}$  and  $\nabla^2 f(x_0) \succeq \mathbf{0}$  (i.e., positive semidefinite).

**Sufficient condition:** Assume  $f(x)$  is 2-order differentiable at  $x_0$ . If  $\nabla f(x_0) = \mathbf{0}$  and  $\nabla^2 f(x_0) \succ \mathbf{0}$  (i.e., positive definite),  $x_0$  is a local min.

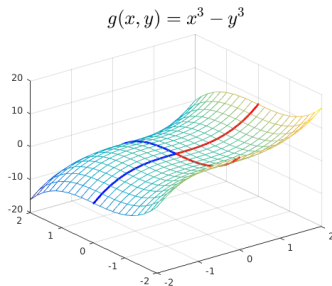
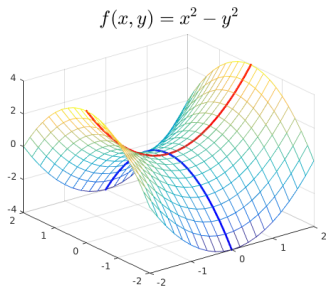
Taylor's:  $f(x_0 + \delta) = f(x_0) + \langle \nabla f(x_0), \delta \rangle + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$ .

- If  $x_0$  is a local min,  $\nabla f(x_0) = \mathbf{0}$  (1st-order condition) and  $f(x_0 + \delta) = f(x_0) + \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$ .
- So  $f(x_0 + \delta) - f(x_0) = \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2) \geq 0$  for all  $\delta$  sufficiently small
- For all  $\delta$  sufficiently small, sign of  $\frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle + o(\|\delta\|_2^2)$  determined by the sign of  $\frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle \implies \frac{1}{2} \langle \delta, \nabla^2 f(x_0) \delta \rangle \geq 0$
- So  $\nabla^2 f(x_0) \succeq \mathbf{0}$ .

# 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  $\nabla^2 f(\mathbf{x}_0) \succeq \mathbf{0}$

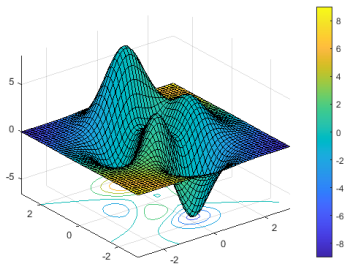


$$\nabla f = \begin{bmatrix} 2x \\ -2y \end{bmatrix}, \nabla^2 f = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$

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

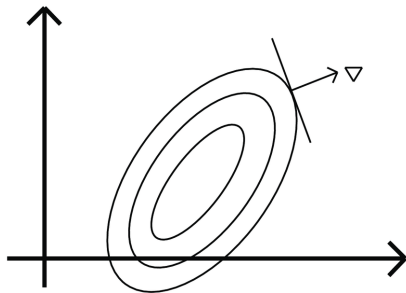


# Contour plot



contour/levelset plot

(Credit: Mathworks)



gradient direction? why?

- [Boyd and Vandenberghe, 2004] Boyd, S. and Vandenberghe, L. (2004). **Convex Optimization**. Cambridge University Press.
- [Coleman, 2012] Coleman, R. (2012). **Calculus on Normed Vector Spaces**. Springer New York.
- [Hiriart-Urruty and Lemaréchal, 2001] Hiriart-Urruty, J.-B. and Lemaréchal, C. (2001). **Fundamentals of Convex Analysis**. Springer Berlin Heidelberg.
- [Munkres, 1997] Munkres, J. R. (1997). **Analysis On Manifolds**. Taylor & Francis Inc.
- [Zorich, 2015] Zorich, V. A. (2015). **Mathematical Analysis I**. Springer Berlin Heidelberg.