Mathematics of Data: From Theory to Computation

Prof. Volkan Cevher volkan.cevher@epfl.ch

Lecture 3: Some basics on optimization

Laboratory for Information and Inference Systems (LIONS) École Polytechnique Fédérale de Lausanne (EPFL)

EE-556 (Fall 2024)

















License Information for Mathematics of Data Slides

▶ This work is released under a <u>Creative Commons License</u> with the following terms:

Attribution

► The licensor permits others to copy, distribute, display, and perform the work. In return, licensees must give the original authors credit.

Non-Commercial

► The licensor permits others to copy, distribute, display, and perform the work. In return, licensees may not use the work for commercial purposes — unless they get the licensor's permission.

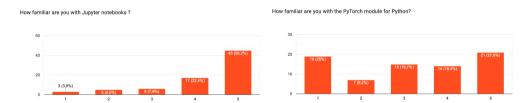
Share Alike

The licensor permits others to distribute derivative works only under a license identical to the one that governs the licensor's work.

► Full Text of the License

Survey responses

- o A majority of respondents are familiar with Python.
 - ► Most are comfortable with Jupyter notebooks.
 - ► There is still some room to learn PyTorch.



Remark:

o Homeworks will be given as Jupyter notebooks.

Outline

- ► This lecture
 - 1. Linear algebra: Norms, matrix norms, dual norms
 - 2. Analysis: Continuity, Lipschitz continuity, differentiation
 - 3. Convexity: Convex sets, convex functions, subdifferentials, L-Lipschitz gradient functions, strong convexity
 - 4. Convergence rates and convergence plots
- Next lecture
 - 1 Gradient descent methods

Vector norms

Definition (Vector norm)

A norm of a vector in \mathbb{R}^p is a function $\|\cdot\|: \mathbb{R}^p \to \mathbb{R}$ such that for all vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^p$ and scalar $\lambda \in \mathbb{R}$

(a) $\|\mathbf{x}\| > 0$ for all $\mathbf{x} \in \mathbb{R}^p$

nonnegativity

(b) $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$

definitiveness homogeniety

(c) $\|\lambda \mathbf{x}\| = |\lambda| \|\mathbf{x}\|$

triangle inequality

- (d) $\|\mathbf{x} + \mathbf{y}\| \le \|\mathbf{x}\| + \|\mathbf{y}\|$
- **Observations:** \circ There is a family of ℓ_q -norms parameterized by $q \in [1, \infty]$;
 - \circ For $\mathbf{x} \in \mathbb{R}^p$, the ℓ_q -norm is defined as $\|\mathbf{x}\|_q := \left(\sum_{i=1}^p |x_i|^q
 ight)^{1/q}$.

Example

- (1) ℓ_2 -norm: $\|\mathbf{x}\|_2 := \sqrt{\sum_{i=1}^p x_i^2}$ (Euclidean norm)
- (2) ℓ_1 -norm: $\|\mathbf{x}\|_1 := \sum_{i=1}^p |x_i|$ (Manhattan norm)
- (3) ℓ_{∞} -norm: $\|\mathbf{x}\|_{\infty} := \max_{i=1,\dots,n} |x_i|$ (Chebyshev norm)

Vector norms contd.

Definition (Quasi-norm)

A quasi-norm satisfies all the norm properties except (d) triangle inequality, which is replaced by $\|\mathbf{x} + \mathbf{y}\| \le c (\|\mathbf{x}\| + \|\mathbf{y}\|)$ for a constant $c \ge 1$.

Definition (Semi(pseudo)-norm)

A semi(pseudo)-norm satisfies all the norm properties except (b) definiteness.

Example

- ▶ The ℓ_q -norm is in fact a quasi norm when $q \in (0,1)$, with $c = 2^{1/q} 1$.
- The total variation norm (TV-norm) defined (in 1D): $\|\mathbf{x}\|_{\mathrm{TV}} := \sum_{i=1}^{p-1} |x_{i+1} x_i|$ is a semi-norm since it fails to satisfy (b); e.g., any $\mathbf{x} = c(1, 1, \ldots, 1)^T$ for $c \neq 0$ will have $\|\mathbf{x}\|_{\mathrm{TV}} = 0$ even though $\mathbf{x} \neq \mathbf{0}$.

Definition (ℓ_0 -"norm")

$$\|\mathbf{x}\|_0 = \lim_{q \to 0} \|\mathbf{x}\|_q^q = |\{i : x_i \neq 0\}|$$

- **Observations:** o The ℓ_0 -"norm" counts the non-zero components of x. Hence, it is not a norm.
 - \circ It does not satisfy the property (c) \Rightarrow it is also neither a quasi- nor a semi-norm.

Vector norms contd.

Norm balls

Radius
$$r$$
 ball in ℓ_q -norm: $\mathcal{B}_q(r) = \{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\|_q \leq r\}$

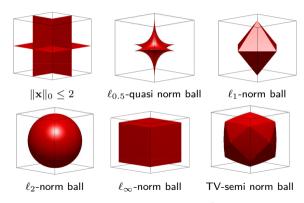


Table: Some norm balls in \mathbb{R}^3

Vector norms contd.

Definition (Dual norm)

Let $\|\cdot\|$ be a norm in \mathbb{R}^p , then the **dual norm** denoted by $\|\cdot\|^*$ is defined:

$$\|\mathbf{x}\|^* = \sup_{\|\mathbf{y}\| \le 1} \mathbf{x}^T \mathbf{y}, \quad \text{for all } \mathbf{x} \in \mathbb{R}^p$$

Observations:

- The dual of the dual norm is the original (primal) norm, i.e., $\|\mathbf{x}\|^{**} = \|\mathbf{x}\|$.
- \circ The dual of $\|\cdot\|_q$ is $\|\cdot\|_p$ where p is such that $\frac{1}{q} + \frac{1}{p} = 1$.
- $\circ \text{ H\"{o}lder's inequality: } |\mathbf{x}^T\mathbf{y}| \leq \|\mathbf{x}\|_q \|\mathbf{y}\|_p \text{, where } p \in [1,+\infty) \text{ and } \tfrac{1}{q} + \tfrac{1}{p} = 1.$
- \circ Cauchy-Schwarz is a special case of Hölder's inequality (q=p=2).

Example

- i) $\|\cdot\|_2$ is dual of $\|\cdot\|_2$ (i.e. $\|\cdot\|_2$ is self-dual): $\sup\{\mathbf{z}^T\mathbf{x} \mid \|\mathbf{x}\|_2 \le 1\} = \|\mathbf{z}\|_2$.
- ii) $\|\cdot\|_1$ is dual of $\|\cdot\|_{\infty}$, (and *vice versa*): $\sup\{\mathbf{z}^T\mathbf{x}\mid \|\mathbf{x}\|_{\infty}\leq 1\}=\|\mathbf{z}\|_1$.

Matrix norms

o Similar to vector norms, matrix norms are a metric over matrices:

Definition (Matrix norm)

A norm of an $n \times p$ matrix is a map $\|\cdot\| : \mathbb{R}^{n \times p} \to \mathbb{R}$ such that for all matrices $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times p}$ and scalar $\lambda \in \mathbb{R}$

- (a) $\|\mathbf{A}\| \geq 0$ for all $\mathbf{A} \in \mathbb{R}^{n \times p}$
- (b) $\|\mathbf{A}\| = 0$ if and only if $\mathbf{A} = \mathbf{0}$
- (c) $\|\lambda \mathbf{A}\| = |\lambda| \|\mathbf{A}\|$
- (d) $\|\mathbf{A} + \mathbf{B}\| < \|\mathbf{A}\| + \|\mathbf{B}\|$

nonnegativity

definitiveness

homogeniety

triangle inequality

Definition (Matrix inner product)

Matrix inner product is defined as follows

$$\langle \mathbf{A}, \mathbf{B} \rangle = \mathsf{trace}\left(\mathbf{A}\mathbf{B}^T\right)$$
 .

o Similar to vector ℓ_p -norms, we have Schatten q-norms for matrices.

Definition (Schatten q-norms)

 $\|\mathbf{A}\|_q := \left(\sum_{i=1}^p (\sigma(\mathbf{A})_i)^q\right)^{1/q}$, where $\sigma(\mathbf{A})_i$ is the i^{th} singular value of \mathbf{A} .

Example (with
$$r = \min\{n, p\}$$
 and $\sigma_i = \sigma(\mathbf{A})_i$)
$$\|\mathbf{A}\|_1^S = \|\mathbf{A}\|_* := \sum_{i=1}^r \sigma_i \qquad \equiv \operatorname{trace}\left(\sqrt{\mathbf{A}^T\mathbf{A}}\right) \quad \text{(Nuclear/trace)}$$

$$\|\mathbf{A}\|_2^S = \|\mathbf{A}\|_F := \sqrt{\sum_{i=1}^r (\sigma_i)^2} \quad \equiv \sqrt{\sum_{i=1}^n \sum_{j=1}^p |a_{ij}|^2} \quad \text{(Frobenius)}$$

$$\|\mathbf{A}\|_\infty^S = \|\mathbf{A}\| \quad := \max_{i=1,\dots,r} \{\sigma_i\} \quad \equiv \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} \quad \text{(Spectral/matrix)}$$

Definition (Operator norm)

The operator norm between ℓ_q and ℓ_r $(1 \le q, r \le \infty)$ of a matrix $\mathbf A$ is defined as

$$\|\mathbf{A}\|_{q\to r} = \sup_{\|\mathbf{x}\|_q \le 1} \|\mathbf{A}\mathbf{x}\|_r$$

Problem

Show that $\|\mathbf{A}\|_{2\to 2} = \|\mathbf{A}\|$ i.e., ℓ_2 to ℓ_2 operator norm is the spectral norm.

Solution

$$\begin{split} \|\mathbf{A}\|_{2\to 2} &= \sup_{\|\mathbf{x}\|_2 \le 1} \|\mathbf{A}\mathbf{x}\|_2 = \sup_{\|\mathbf{x}\|_2 \le 1} \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{x}\|_2 \quad \text{(using SVD of } \mathbf{A} \text{)} \\ &= \sup_{\|\mathbf{x}\|_2 \le 1} \|\mathbf{\Sigma}\mathbf{V}^T\mathbf{x}\|_2 \quad \text{(rotational invariance of } \|\cdot\|_2 \text{)} \\ &= \sup_{\|\mathbf{z}\|_2 \le 1} \|\mathbf{\Sigma}\mathbf{z}\|_2 \quad \text{(letting } \mathbf{V}^T\mathbf{x} = \mathbf{z} \text{)} \\ &= \sup_{\|\mathbf{z}\|_2 \le 1} \sqrt{\sum_{i=1}^{\min(n,p)} \sigma_i^2 z_i^2} = \sigma_{\max} = \|\mathbf{A}\| \end{split}$$

Other examples

▶ The $\|\mathbf{A}\|_{\infty\to\infty}$ (norm induced by ℓ_{∞} -norm) also denoted $\|\mathbf{A}\|_{\infty}$, is the max-row-sum norm:

$$\|\mathbf{A}\|_{\infty \to \infty} := \sup\{\|\mathbf{A}\mathbf{x}\|_{\infty} \mid \|\mathbf{x}\|_{\infty} \le 1\} = \max_{i=1,\dots,n} \sum_{j=1}^{p} |a_{ij}|.$$

▶ The $\|\mathbf{A}\|_{1\to 1}$ (norm induced by ℓ_1 -norm) also denoted $\|\mathbf{A}\|_1$, is the max-column-sum norm:

$$\|\mathbf{A}\|_{1\to 1} := \sup\{\|\mathbf{A}\mathbf{x}\|_1 \mid \|\mathbf{x}\|_1 \le 1\} = \max_{i=1,\dots,p} \sum_{j=1}^n |a_{ij}|.$$

Matrix & vector norm analogy

Vectors	$\ \mathbf{x}\ _1$	$\ \mathbf{x}\ _2$	$\ \mathbf{x}\ _{\infty}$
Matrices	$\ \mathbf{X}\ _*$	$\ \mathbf{X}\ _F$	$\ \mathbf{X}\ $

Definition (Dual of a matrix)

The dual norm of $\mathbf{A} \in \mathbb{R}^{n \times p}$ is defined as

$$\|\mathbf{A}\|^* = \sup \left\{ \operatorname{trace} \left(\mathbf{A}^T \mathbf{X} \right) \mid \|\mathbf{X}\| \le 1 \right\}.$$

Matrix & vector dual norm analogy

Vector primal norm	$\ \mathbf{x}\ _1$	$\ \mathbf{x}\ _2$	$\ \mathbf{x}\ _{\infty}$
Vector dual norm	$\ \mathbf{x}\ _{\infty}$	$\ \mathbf{x}\ _2$	$\ \mathbf{x}\ _1$
Matrix primal norm	$\ \mathbf{X}\ _*$	$\ \mathbf{X}\ _F$	$\ \mathbf{X}\ $
Matrix dual norm	$\ \mathbf{X}\ $	$\ \mathbf{X}\ _F$	$\ \mathbf{X}\ _*$

Matrix definitions contd.

Definition (Positive semidefinite & positive definite matrices)

A symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is positive semidefinite (denoted $\mathbf{A} \succeq 0$) if $\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ for all $\mathbf{x} \neq \mathbf{0}$; while it is positive definite (denoted $\mathbf{A} \succ 0$) if $\mathbf{x}^T \mathbf{A} \mathbf{x} > 0$.

Observations:

- \circ **A** $\succeq 0$ iff all its eigenvalues are **nonnegative** i.e. $\lambda_{\min}(\mathbf{A}) \geq 0$.
- \circ Similarly, $\mathbf{A} \succ 0$ iff all its eigenvalues are **positive** i.e. $\lambda_{\min}(\mathbf{A}) > 0$.
- \circ **A** is negative semidefinite if $-\mathbf{A} \succeq 0$; while **A** is negative definite if $-\mathbf{A} \succ 0$.
- \circ Semidefinite ordering of two symmetric matrices, A and B: $A \succeq B$ if $A B \succeq 0$.

Example (Matrix inequalities)

- 1. If $\mathbf{A} \succeq 0$ and $\mathbf{B} \succeq 0$, then $\mathbf{A} + \mathbf{B} \succeq 0$
- 2. If $A \succeq B$ and $C \succeq D$, then $A + C \succeq B + D$
- 3. If $\mathbf{B} \leq 0$ then $\mathbf{A} + \mathbf{B} \leq \mathbf{A}$
- 4. If $\mathbf{A} \succeq 0$ and $\alpha \geq 0$, then $\alpha \mathbf{A} \succeq 0$
- 5. If $\mathbf{A} \succ 0$, then $\mathbf{A}^2 \succ 0$
- 6. If $\mathbf{A} \succ 0$, then $\mathbf{A}^{-1} \succ 0$

Continuity in functions

Definition (Continuity)

Let $f:\mathcal{Q}\to\mathbb{R}$ where $\mathcal{Q}\subseteq\mathbb{R}^p$. Then, f is a continuous function over its domain \mathcal{Q} if and only if

$$\lim_{\mathbf{x}\to\mathbf{y}} f(\mathbf{x}) = f(\mathbf{y}), \quad \forall \mathbf{y} \in \mathcal{Q},$$

i.e., the limit of f—as $\mathbf x$ approaches $\mathbf y$ —exists and is equal to $f(\mathbf y)$.

Definition (Class of continuous functions)

We denote the class of continuous functions f over the domain \mathcal{Q} as $f \in \mathcal{C}(\mathcal{Q})$.

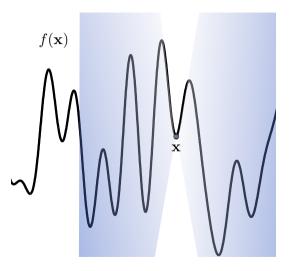
Definition (Lipschitz continuity)

Let $f:\mathcal{Q}\to\mathbb{R}$ where $\mathcal{Q}\subseteq\mathbb{R}^p$. Then, f is called Lipschitz continuous if there exists a constant value $K\geq 0$ such that the following holds

$$|f(\mathbf{y}) - f(\mathbf{x})| \le K ||\mathbf{y} - \mathbf{x}||_2, \quad \forall \mathbf{x}, \ \mathbf{y} \in \mathcal{Q}.$$

Observation: o "Small" changes in the input result into "small" changes in the function values.

Continuity in functions



Differentiability in functions

Definition (Differentiability)

Let $\mathcal{Q} \subseteq \mathbb{R}^p$. A function $f: \mathcal{Q} \to \mathbb{R}$ is said to be k-times continuously differentiable on \mathcal{Q} if all its partial derivatives up to k-th order exist and are continuous over \mathcal{Q} . Notation: $f \in \mathcal{C}^k(\mathcal{Q})$.

o A key quantity is the gradient of the function $f: \mathcal{Q} \to \mathbb{R}$, which we denote as ∇f (e_i is the i-th unit vector):

$$\nabla f(\mathbf{x}) := \sum_{i=1}^{p} \frac{\partial f}{\partial x_i} \mathbf{e}_i = \left[\frac{\partial f}{\partial x_1}, \cdots, \frac{\partial f}{\partial x_p} \right]^T.$$

 $\circ \text{ For } k=2 \text{, we dub } \nabla^2 f \text{ as the } \mathbf{Hessian} \text{ of } f \text{, i.e., } \left[\nabla^2 f \right]_{i,j} := \frac{\partial^2 f}{\partial x_i \partial x_j}.$

Gradients as linear approximations

A "Taylor" way of thinking about gradients:

Let $\mathcal{Q}\subseteq\mathbb{R}^p$. If $f\in\mathcal{C}^1(\mathcal{Q})$, then $\mathbf{u}\mapsto\langle\nabla f(\mathbf{x}),\mathbf{u}\rangle$ is the *unique* linear function from \mathcal{Q} to \mathbb{R} such that

$$\lim_{\mathbf{u} \rightrightarrows 0} \frac{|f(\mathbf{x} + \mathbf{u}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{u} \rangle|}{\|\mathbf{u}\|} \to 0$$

Example

The gradient of $f: \mathbf{x} \mapsto \|\mathbf{x}\|_2^2$ is

$$\nabla f(\mathbf{x}) = 2\mathbf{x}$$

Proof:

o To apply the Taylor way of thinking, we consider the following quantity:

$$\begin{split} f(\mathbf{x} + \mathbf{u}) - f(\mathbf{x}) &= \|\mathbf{x} + \mathbf{u}\|_2^2 - \|\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2 + 2\langle \mathbf{x}, \mathbf{u} \rangle + \|\mathbf{u}\|_2^2 - \|\mathbf{x}\|_2^2 \\ &= 2\langle \mathbf{x}, \mathbf{u} \rangle + \|\mathbf{u}\|_2^2 \\ &= \langle 2\mathbf{x}, \mathbf{u} \rangle + o(\|\mathbf{u}\|_2). \end{split}$$

o Since the linear map is unique, we get that the gradient is $\nabla f(\mathbf{x}) = 2\mathbf{x}$.

To be or not to be differentiable

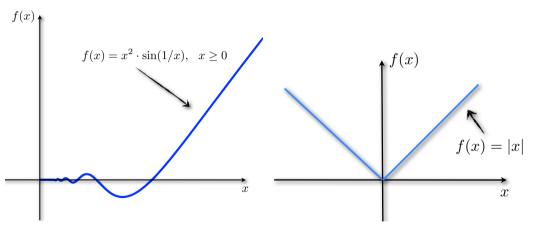


Figure: (Left panel) ∞ -times continuously differentiable function in $\mathbb R$. (Right panel) Non-differentiable f(x) = |x| in $\mathbb R$.

Gradients of vector valued functions

Jacobian

When $f:\mathbb{R}^n
ightrightarrows \mathbb{R}^d$ is a vector valued function, the following d imes n matrix ${f J}$ of partial derivatives

$$\left[\mathbf{J}_f(\mathbf{x})\right]_{i,j} := \frac{\partial f_i}{\partial x_j}(\mathbf{x})$$

is called the Jacobian of f at x.

Observations: \circ The Jacobian is the transpose of the gradient, when f is real valued.

o Thinking in terms of Jacobians is really helpful when we need to use the chain rule.

Chain Rule via Jacobians

Let \circ denote the functional composition: $g \circ f := g(f(\mathbf{x}))$. If $g \circ f$ is differentiable at \mathbf{x} , then the following holds

$$\mathbf{J}_{g \circ f}(\mathbf{x}) = \mathbf{J}_g(f(\mathbf{x}))\mathbf{J}_f(\mathbf{x}).$$

Hence, the chain rule, which is helpful in differentiating function compositions, can be related to a simple product of Jacobian matrices.

Example: Quadratic loss

Example

The gradient of the function $h: \mathbf{x} \mapsto \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$ is given by the following expression:

$$\nabla h(\mathbf{x}) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b}).$$

Proof: o

- We apply the chain rule:
 - ▶ The Jacobian of the affine function $f : \mathbf{x} \mapsto \mathbf{A}\mathbf{x} \mathbf{b}$ is $\mathbf{J}_f(\mathbf{x}) = \mathbf{A}$.
 - ▶ The gradient of $g: \mathbf{x} \mapsto \|\mathbf{x}\|_2^2$ is $\nabla g(\mathbf{x}) = 2\mathbf{x} \Rightarrow \mathbf{J}_g(\mathbf{x}) = 2\mathbf{x}^T$.
 - ▶ Using the chain rule on the composition $h = g \circ f$:

$$\begin{aligned} \mathbf{J}_{g \circ f}(\mathbf{x}) &= \mathbf{J}_g(f(\mathbf{x})) \mathbf{J}_f(\mathbf{x}) \\ &= \mathbf{J}_g(\mathbf{A}\mathbf{x} - \mathbf{b}) \mathbf{J}_f(\mathbf{x}) \\ &= 2(\mathbf{A}\mathbf{x} - \mathbf{b})^T \mathbf{A}. \end{aligned}$$

 \circ Since h is real valued, the Jacobian is a row vector, we obtain the gradient by transposing.

Example: Logistic loss

Example

The gradient of the logistic loss $f(\mathbf{x}) = \log(1 + \exp(-b(\mathbf{a}^T\mathbf{x})))$ is given by the following expression:

$$\nabla f(\mathbf{x}) = -b \frac{\exp(-b(\mathbf{a}^T \mathbf{x}))}{1 + \exp(-b(\mathbf{a}^T \mathbf{x}))} \mathbf{a}.$$

Proof:

- \circ f is a composition of the following functions:
 - $lackbox{h}(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$, whose Jacobian is $\mathbf{J}_h(\mathbf{x}) = \mathbf{a}^T$
 - $lacksquare g(u) = \log(1 + \exp(-bu))$, whose " 1×1 Jacobian" is $\mathbf{J}_g(u) = -b \frac{\exp(-bu)}{1 + \exp(-bu)}$
 - ► By the chain rule:

$$\mathbf{J}_f(\mathbf{x}) = \mathbf{J}_g(h(\mathbf{x})) \cdot \mathbf{J}_h(\mathbf{x})$$
$$= -b \frac{\exp(-b(\mathbf{a}^T \mathbf{x}))}{1 + \exp(-b(\mathbf{a}^T \mathbf{x}))} \mathbf{a}^T$$

 \circ The gradient is simply the transpose of $\mathbf{J}_f(\mathbf{x})$.

Use Jacobians !

With Jacobians, differentiating function compositions is a direct mechanical process.

A more complicated example here and another one at the advanced material!

Example

The gradient of $f: \mathbf{x} \mapsto w_2^T \sigma(\mathbf{W}_1 \mathbf{x} + \boldsymbol{\mu})$ is given by the following expression:

$$\nabla f(\mathbf{x}) = \mathbf{J}_f(\mathbf{x})^T = \mathbf{W}_1^T(\sigma'(\mathbf{W}_1\mathbf{x} + \boldsymbol{\mu}) \odot \boldsymbol{w}_2),$$

where σ is a non-linear function that applies to each coordinate, and \odot denotes the component wise product.

Proof:

- \circ We use the fact that f is a composition of the following functions:
 - $h(\mathbf{x}) = \mathbf{W}_1 \mathbf{x} + \mu$, whose Jacobian is $\mathbf{J}_h(\mathbf{x}) = \mathbf{W}_1$.

- $lackbox{f k}({f x})={m w}_2^T{f x}$ whose Jacobian is ${f J}_k({f x})={m w}_2^T.$
- By the chain rule, we have that

$$\begin{aligned} \mathbf{J}_f(\mathbf{x}) &= \mathbf{J}_k(g(h(\mathbf{x}))) \cdot \mathbf{J}_g(h(\mathbf{x})) \cdot \mathbf{J}_h(\mathbf{x}) \\ &= w_2^T \cdot \mathsf{diag}(\sigma'([\mathbf{W}_1\mathbf{x} + \boldsymbol{\mu}]_1), \dots, \sigma'([\mathbf{W}_1\mathbf{x} + \boldsymbol{\mu}]_n)) \cdot \mathbf{W}_1. \end{aligned}$$

o Simply transpose the Jacobian to get the gradient and use o to replace the diagonal matrix.

Some reminders on sets

Definition (Closed set)

A set is closed if it contains all its limit points.

Definition (Open set)

A set is open if its complement is closed.

Definition (Closure of a set)

Let $Q \subseteq \mathbb{R}^p$ be a given open set, i.e., it contains a neighborhood of all its points. Then, the closure of Q, denoted as $\operatorname{cl}(Q)$, is the smallest closed set in \mathbb{R}^p that includes Q.



Figure: (Left panel) Closed set Q. (Middle panel) Open set Q and its closure cl(Q) (Right panel).

Convexity of sets

Definition

 $ightharpoonup \mathcal{Q} \subseteq \mathbb{R}^p$ is a convex set if

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

 $ightharpoonup \mathcal{Q} \subseteq \mathbb{R}^p$ is a *strictly* convex set if

$$\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q} \quad \forall \alpha \in (0, 1), \quad \alpha \mathbf{x}_1 + (1 - \alpha) \mathbf{x}_2 \in \mathsf{interior}(\mathcal{Q}).$$

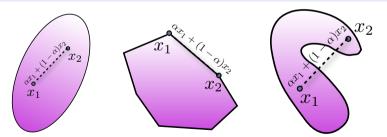


Figure: (Left) Strictly convex (Middle) Convex (Right) Non-convex

Definition

Let \mathcal{Q} be a convex set in \mathbb{R}^p . A function $f \colon \mathcal{Q} \to \mathbb{R}$ is called *convex* if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) \le \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2), \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

ightharpoonup f is called concave, if -f is convex.

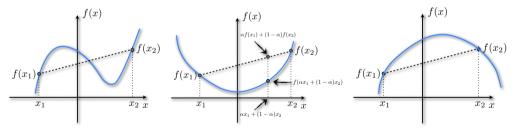


Figure: (Left) Non-convex (Middle) Convex (Right) Concave

Definition

Let $\mathcal Q$ be a convex set in $\mathbb R^p$. A function $f\colon \mathcal Q\to\mathbb R$ is called *convex* if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) \le \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2), \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

Question: \circ Can we extend f from $\mathcal Q$ to $\mathbb R^p$ preserving convexity?

Definition

Let $\mathcal Q$ be a convex set in $\mathbb R^p$. A function $f\colon \mathcal Q\to\mathbb R$ is called *convex* if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) \le \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2), \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

Question: \circ Can we extend f from $\mathcal Q$ to $\mathbb R^p$ preserving convexity?

Definition (Extended real-valued convex functions)

$$f(\mathbf{x}) := \left\{ egin{array}{ll} f(\mathbf{x}) & ext{if } \mathbf{x} \in \mathcal{Q} \\ +\infty & ext{if otherwise} \end{array}
ight.$$

Recall, dom(f) = Q. If $Q \neq \mathbb{R}^p$, extended f is never continuous, but it is l.s.c.

Definition

Let $\mathcal Q$ be a convex set in $\mathbb R^p$. A function $f\colon \mathcal Q\to\mathbb R$ is called *convex* if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) \le \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2), \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

Proposition

Every ℓ_q -norm $\|\cdot\|_q$ $(q\geq 1)$ in \mathbb{R}^p is convex.

Proof:

Definition

Let $\mathcal Q$ be a convex set in $\mathbb R^p$. A function $f\colon \mathcal Q\to\mathbb R$ is called *convex* if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) \le \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2), \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

Proposition

Every ℓ_q -norm $\|\cdot\|_q$ $(q \ge 1)$ in \mathbb{R}^p is convex.

Proof: o Proof by intimidation.

Definition

Let $\mathcal Q$ be a convex set in $\mathbb R^p$. A function $f\colon \mathcal Q\to\mathbb R$ is called *convex* if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) \le \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2), \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

Proposition

Every ℓ_q -norm $\|\cdot\|_q$ $(q \ge 1)$ in \mathbb{R}^p is convex.

Proof: • Kidding! By triangle inequality and homogeneity of the norm:

$$\|\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2\|_q \le \|\alpha \mathbf{x}_1\|_q + \|(1 - \alpha)\mathbf{x}_2\|_q = \alpha\|\mathbf{x}_1\|_q + (1 - \alpha)\|\mathbf{x}_2\|_q, \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \forall \alpha \in [0, 1].$$

Definition

Let $\mathcal Q$ be a convex set in $\mathbb R^p$. A function $f\colon \mathcal Q\to\mathbb R$ is called *convex* if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) \le \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2), \qquad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

Example

Function	Example	Attributes
ℓ_q vector norms, $q \geq 1$	$\ \mathbf{x}\ _{2}, \ \mathbf{x}\ _{1}, \ \mathbf{x}\ _{\infty}$	convex
ℓ_q matrix norms, $q \geq 1$	$\ \mathbf{X}\ _* = \sum_{i=1}^{rank(\mathbf{X})} \sigma_i$	convex
Square root function	\sqrt{x}	concave
Max of convex functions	$\max_i f_i(x)$, f_i convex	convex
Min of concave functions	$\min_i f_i(x)$, f_i concave	concave
Sum of convex functions	$\sum_{i=1}^n f_i, f_i$ convex	convex
Logarithmic functions	$\log\left(det(\mathbf{X}) ight)$	concave, assumes $\mathbf{X}\succ 0$
Affine/linear functions	$\sum_{i=1}^{n} X_{ii}$	both convex and concave
Eigenvalue functions	$\lambda_{\max}(\mathbf{X})$	convex, assumes $\mathbf{X} = \mathbf{X}^T$

Revisiting: Alternative definitions of function convexity II [2]

Recall, the epigraph of $f \colon \mathcal{Q} \to \mathbb{R} \cup \{+\infty\}$ is

$$\mathrm{epi}(f) = \left\{ (\mathbf{x}, u) \in \mathcal{Q} \times \mathbb{R} \colon f(\mathbf{x}) \le u \right\}.$$

Definition

A function $f: \mathcal{Q} \to \mathbb{R} \cup \{+\infty\}$ is convex if its epigraph is a convex set.

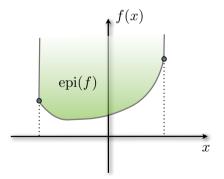
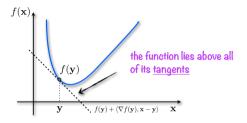


Figure: Epigraph — the region in green above graph f.

Revisiting: Alternative definition of function convexity III [2]



Definition

Let $\mathcal Q$ is a convex set in $\mathbb R^p$. A function $f\in\mathcal C^1(\mathcal Q)$ is called convex on $\mathcal Q$ if for any $\mathbf x,\ \mathbf y\in\mathcal Q$:

$$f(\mathbf{x}) \ge f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \ \mathbf{x} - \mathbf{y} \rangle.$$

Definition

A function $f \in \mathcal{C}^1(\mathcal{Q})$ is called convex on \mathcal{Q} if for any $\mathbf{x}, \ \mathbf{y} \in \mathcal{Q}$:

$$\langle \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}), \ \mathbf{y} - \mathbf{x} \rangle \ge 0.$$

*That is, if its gradient is a monotone operator.



Revisiting: Alternative definition of function convexity IV [2]

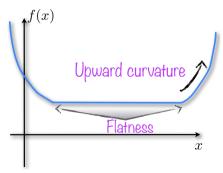
Definition

Let $\mathcal Q$ is a convex set in $\mathbb R^p$. A function $f\in\mathcal C^2(\mathcal Q)$ is called convex on $\mathcal Q$ if for any $\mathbf x\in\mathcal Q$:

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

Remarks:

- \circ Geometrical interpretation: the graph of f has zero or positive (upward) curvature.
- \circ However, this does not exclude flatness of f.



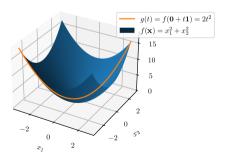
Revisiting: Alternative definition of function convexity V [2]

Definition

Let $\mathcal Q$ is a convex set in $\mathbb R^p$. A function $f\in\mathcal C^2(\mathcal Q)$ is called convex on $\mathcal Q$ if for any $\mathbf x\in\mathcal Q$, $\mathbf v\in\mathbb R^p$, the function $g(t)=f(\mathbf x+t\mathbf v)$ is convex on its domain $\{t|\mathbf x+t\mathbf v\in\mathcal Q\}$.

Remarks:

- \circ This approach allows us to check the convexity long 1-dimensional lines.
- o This concept generalizes to self-concordant functions (advanced material).



Strict convexity

Definition

A function $f\colon \mathcal{Q}\to \mathbb{R}\cup\{+\infty\}$ is called *strictly convex* on $\mathcal Q$ if

$$f(\alpha \mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2) < \alpha f(\mathbf{x}_1) + (1 - \alpha)f(\mathbf{x}_2) \quad \forall \mathbf{x}_1 \ \mathbf{x}_2 \in \mathcal{Q}, \quad \forall \alpha \in (0, 1).$$

Theorem

If $Q \subset \mathbb{R}^p$ is a convex set and $f : \mathbb{R}^p \to (-\infty, +\infty]$ is a proper and strictly convex function, then there exist at most one minimizer of f over Q.

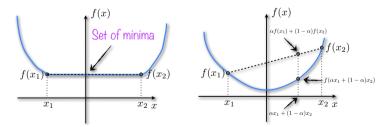


Figure: (Left panel) Convex function. (Right panel) Strictly convex function.

Subdifferentials and (sub)gradients in convex functions

Definition

Let $f:\mathcal{Q}\to\mathbb{R}\cup\{+\infty\}$ be a convex function. The subdifferential of f at a point $\mathbf{x}\in\mathcal{Q}$ is defined by the set:

$$\partial f(\mathbf{x}) = \left\{ \mathbf{v} \in \mathbb{R}^p \ : \ f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{v}, \ \mathbf{y} - \mathbf{x} \rangle \text{ for all } \mathbf{y} \in \mathcal{Q} \right\}.$$

Each element \mathbf{v} of $\partial f(\mathbf{x})$ is called *subgradient* of f at \mathbf{x} .

Definition

Let $f: \mathcal{Q} \to \mathbb{R} \cup \{+\infty\}$ be a differentiable convex function. Then, the subdifferential of f at a point $\mathbf{x} \in \mathcal{Q}$ contains only the gradient, i.e., $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}.$

Remark:

 \circ Subdifferential generalizes abla to nondifferentiable functions

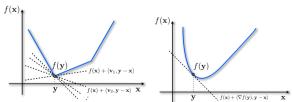


Figure: (Left) Non-differentiability at point y. (Right) Gradient as a subdifferential with a singleton entry.

Generalized subdifferentials for nonconvex functions

Definition

Let $f: \mathcal{Q} \to \mathbb{R} \cup \{+\infty\}$ be a locally Lipschitz function. The Clarke subdifferential of f at a point $\mathbf{x} \in \mathcal{Q}$ is defined by the set:

$$\partial_C f(\mathbf{x}) = \operatorname{conv} \left(\left\{ \mathbf{v} \in \mathbb{R}^p : \begin{array}{l} \exists \mathbf{x}^k o \mathbf{x}, \nabla f\left(\mathbf{x}^k\right) \text{ exists,} \\ \nabla f\left(\mathbf{x}^k\right) o \mathbf{v} \end{array} \right\} \right).$$

Remarks:

- \circ For convex functions, the Clarke subdifferential reduces to subdifferential.
- o If \mathbf{x}^* is a local minimum of f, then $\mathbf{0} \in \partial_C f(\mathbf{x}^*)$.

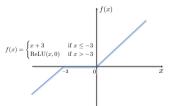


Figure: The Clarke subdifferential at -3 and 0: $\partial_C f(-3) \stackrel{!}{=} \partial_C f(0) = [0,1]$. Non-subdifferentiability at -3 and 0.

Heads up: Be careful with automatic differentiation!

Example (Simple)

The gradient of the function $f: x \mapsto \mathsf{ReLU}(x) - \mathsf{ReLU}(-x) = x$ at 0 is given by g(0) = 1.

Remark:

- Subdifferentials are tricky business!
- Automatic differentiation can be wrong [3]!
- o We will revisit when we discuss the Moreau-Rockafellar's decomposition theorem.

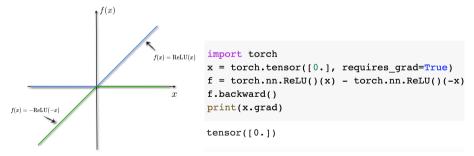


Figure: (Left panel) ReLU function. (Right panel) Calculation of q(0) in PyTorch.

L-Lipschitz gradient class of functions

Definition (*L*-Lipschitz gradient convex functions)

Let $f: \mathcal{Q} \to \mathbb{R}$ be differentiable and convex, i.e., $f \in \mathcal{F}^1(\mathcal{Q})$. Then, f has a Lipschitz gradient if there exists L > 0 (the Lipschitz constant) such that $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \le L\|\mathbf{x} - \mathbf{y}\|_2$, $\forall \mathbf{x}, \mathbf{y} \in \mathcal{Q}$.

Proposition (*L*-Lipschitz gradient convex functions)

 $f \in \mathcal{F}^1(\mathcal{Q})$ has L-Lipschitz gradient if and only if the following function is convex:

$$h(\mathbf{x}) = \frac{L}{2} \|\mathbf{x}\|_2^2 - f(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{Q}.$$

Definition (Class of 2-nd order Lipschitz functions)

The class of twice continuously differentiable functions f on $\mathcal Q$ with Lipschitz continuous Hessian is denoted as $\mathcal F_L^{2,2}(\mathcal Q)$ (with $2\to 2$ denoting the spectral norm)

$$\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\|_{2\to 2} \le L \|\mathbf{x} - \mathbf{y}\|_2, \quad \forall \mathbf{x}, \mathbf{y} \in Q,$$

Remark: o $\mathcal{F}_{L}^{l,m}$: functions that are l-times differentiable with m-th order Lipschitz property.

Example: Logistic regression

Problem (Logistic regression)

Given a sample vector $\mathbf{a}_i \in \mathbb{R}^p$ and a binary class label $b_i \in \{-1, +1\}$ (i = 1, ..., n), we define the conditional probability of b_i given \mathbf{a}_i as:

$$\mathbb{P}(b_i|\mathbf{a}_i,\mathbf{x}^{\natural},\mu) \propto 1/(1+e^{-b_i(\langle\mathbf{x}^{\natural},\mathbf{a}_i\rangle+\mu)}),$$

where $\mathbf{x}^{\natural} \in \mathbb{R}^p$ is some true weight vector, $\mu \in \mathbb{R}$ is called the intercept. How to estimate \mathbf{x}^{\natural} given the sample vectors, the binary labels, and μ ?

Optimization formulation

$$\min_{\mathbf{x} \in \mathbb{R}^p} \underbrace{\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-b_i(\mathbf{a}_i^T \mathbf{x} + \mu)))}_{f(\mathbf{x})}$$

Structural properties

Let $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]^T$ (design matrix), then $f \in \mathcal{F}_L^{2,1}$, with $\underline{\mathbf{L}} = \frac{1}{4} \|\mathbf{A}^T \mathbf{A}\|$

μ -strongly convex functions

Definition

A function $f:\mathcal{Q}\to\mathbb{R}\cup\{+\infty\}$, $\mathcal{Q}\subseteq\mathbb{R}^p$ is called μ -strongly convex on its domain if and only if for any \mathbf{x} , $\mathbf{y}\in\mathcal{Q}$ and $\alpha\in[0,1]$ we have:

$$f(\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}) \le \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) - \frac{\mu}{2}\alpha(1 - \alpha)\|\mathbf{x} - \mathbf{y}\|_2^2.$$

The constant μ is called the convexity parameter of function f.

- ▶ The class of k-differentiable μ -strongly functions is denoted as $\mathcal{F}^k_{\mu}(\mathcal{Q})$.
- ► Strong convexity ⇒ strict convexity, BUT strict convexity ⇒ strong convexity

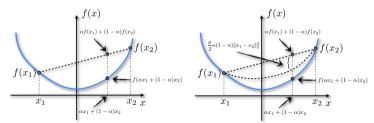


Figure: (Left) Convex (Right) Strongly convex

Alternative: μ -strongly convex functions

Definition

A convex function $f:\mathcal{Q}\to\mathbb{R}$ is said to be $\mu\text{-strongly convex}$ if

$$h(\mathbf{x}) = f(\mathbf{x}) - \frac{\mu}{2} \|\mathbf{x}\|_2^2$$

is convex, where μ is called the strong convexity parameter.

- ▶ The class of k-differentiable μ -strongly functions is denoted as $\mathcal{F}^k_{\mu}(\mathcal{Q})$.
- Non-smooth functions can be μ -strongly convex: e.g., $f(\mathbf{x}) = \|\mathbf{x}\|_1 + \frac{\mu}{2} \|\mathbf{x}\|_2^2$.

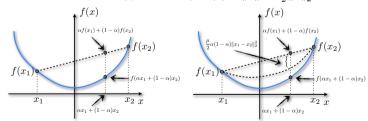


Figure: (Left) Convex (Right) Strongly convex

Lemma

Let $f: \mathcal{Q} \to \mathbb{R}$, $\mathcal{Q} \subseteq \mathbb{R}^p$ be a twice differentiable convex function, i.e., $f \in \mathcal{F}^2(\mathcal{Q})$. Then, f is μ -strongly convex function if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}, \ \forall \mathbf{x} \in \mathbb{R}^p.$$

Lemma

Let $f: \mathcal{Q} \to \mathbb{R}$, $\mathcal{Q} \subseteq \mathbb{R}^p$ be a twice differentiable convex function, i.e., $f \in \mathcal{F}^2(\mathcal{Q})$. Then, f is μ -strongly convex function if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}, \ \forall \mathbf{x} \in \mathbb{R}^p.$$

Example (Toy example)

Consider the quadratic function $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$. Then, f is a μ -strongly convex since $\nabla^2 f(\mathbf{x}) = \mathbf{I} \implies \mu = 1$.

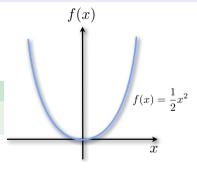


Figure: Toy example for μ -strongly convex functions.

Lemma

Let $f: \mathcal{Q} \to \mathbb{R}, \mathcal{Q} \subseteq \mathbb{R}^p$ be a twice differentiable convex function, i.e., $f \in \mathcal{F}^2(\mathcal{Q})$. Then, f is μ -strongly convex function if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}, \ \forall \mathbf{x} \in \mathbb{R}^p.$$

Example (Overdetermined least squares)

Consider an overdetermined linear system of equations $\mathbf{b} = \mathbf{A}\mathbf{x}^{\natural} + \mathbf{w}$ where $\mathbf{A} \in \mathbb{R}^{n \times p}$ is a full column-rank matrix and \mathbf{x}^{\natural} is unknown. Assume that $\mathbf{A}^T\mathbf{A} \succeq \rho\mathbf{I}, \rho > 0$ and let $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2$. Then, f is a μ -strongly convex function, i.e., $f \in \mathcal{F}^2_{\mu}(\mathbb{R}^p)$ since:

$$\nabla^2 f(\mathbf{x}) = \mathbf{A}^T \mathbf{A}$$
 where $\mathbf{A}^T \mathbf{A} \succeq \rho \mathbf{I} =: \mu \mathbf{I}$.

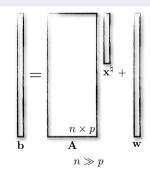


Figure: Overdetermined system of linear equations.

Lemma

Let $f: \mathcal{Q} \to \mathbb{R}$, $\mathcal{Q} \subseteq \mathbb{R}^p$ be a twice differentiable convex function, i.e., $f \in \mathcal{F}^2(\mathcal{Q})$. Then, f is μ -strongly convex function if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}, \ \forall \mathbf{x} \in \mathbb{R}^p.$$

Example (Trivial)

Any linear function $f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} + \beta \in \mathcal{F}^1_\mu(\mathbb{R}^p)$ for $\mu = 0$ since

$$\nabla f(\mathbf{x}) = \mathbf{c}$$
 and $\nabla^2 f(\mathbf{x}) = \mathbf{0}$.

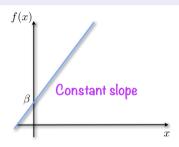


Figure: Counterexample for μ -strongly convex functions.

Lemma

Let $f: \mathcal{Q} \to \mathbb{R}, \mathcal{Q} \subseteq \mathbb{R}^p$ be a twice differentiable convex function, i.e., $f \in \mathcal{F}^2(\mathcal{Q})$. Then, f is μ -strongly convex function if and only if

$$\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}, \ \forall \mathbf{x} \in \mathbb{R}^p.$$

Lemma

A continuously differentiable function f belongs to $\mathcal{F}^1_{\mu}(\mathcal{Q})$ if there exists a constant $\mu > 0$ such that for any $\mathbf{x}, \mathbf{y} \in \mathcal{Q}$, we have:

$$f(\mathbf{y}) \ge f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} ||\mathbf{y} - \mathbf{x}||_2^2$$

Lemma

Let f be continuously differentiable. The following condition, holding for all $\mathbf{x}, \mathbf{y} \in \mathcal{Q} \subseteq \mathbb{R}^p$, is equivalent to inclusion that f is μ -strongly convex function:

$$\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \ge \mu \|\mathbf{x} - \mathbf{y}\|_2^2.$$

L-smooth, μ -strongly convex functions

Definition

Let $f:\mathcal{Q}\to\mathbb{R}, \mathcal{Q}\subseteq\mathbb{R}^p$ be a continuously differentiable function. Then, f is both μ -strongly and L-smooth convex function if for any $\mathbf{x},\mathbf{y}\in\mathcal{Q}$, we have:

$$\frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \le f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \le \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

for constants $0 < \mu \le L$. We denote that $f \in \mathcal{F}^{1,1}_{\mu,L}(\mathcal{Q})$. If f is twice differentiable, an equivalent condition is

$$\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}.$$

L-smooth, μ -strongly convex functions

Definition

Let $f:\mathcal{Q}\to\mathbb{R},\mathcal{Q}\subseteq\mathbb{R}^p$ be a continuously differentiable function. Then, f is both μ -strongly and L-smooth convex function if for any $\mathbf{x},\mathbf{y}\in\mathcal{Q}$, we have:

$$\frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \le f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \le \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

for constants $0 < \mu \le L$. We denote that $f \in \mathcal{F}_{\mu,L}^{1,1}(\mathcal{Q})$. If f is twice differentiable, an equivalent condition is

$$\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}.$$

Example

Consider an linear system of equations $\mathbf{b} = \mathbf{A}\mathbf{x}^{\natural}$ where $\mu \mathbf{I} \leq \mathbf{A}^T \mathbf{A} \leq L \mathbf{I}$. Let $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2$. Then, f is both μ -strongly convex and L-smooth function, i.e., $f \in \mathcal{F}_{\mu,L}^{2,1}(\mathbb{R}^p)$ since:

$$\nabla^2 f(\mathbf{x}) = \mathbf{A}^T \mathbf{A}$$
 where $\mu \mathbf{I} \preceq \mathbf{A}^T \mathbf{A} \preceq L \mathbf{I}$.

L-smooth, μ -strongly convex functions

Definition

Let $f: \mathcal{Q} \to \mathbb{R}, \mathcal{Q} \subseteq \mathbb{R}^p$ be a continuously differentiable function. Then, f is both μ -strongly and L-smooth convex function if for any $\mathbf{x}, \mathbf{y} \in \mathcal{Q}$, we have:

$$\frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 \le f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \le \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

for constants $0 < \mu \le L$. We denote that $f \in \mathcal{F}^{1,1}_{\mu,L}(\mathcal{Q})$. If f is twice differentiable, an equivalent condition is

$$\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}.$$

Observations:

- \circ Both μ and L show up in convergence rate characterization of algorithms
- \circ Unfortunately, μ,L are usually not known a priori...
- o When they are known, they can help significantly (even in stopping algorithms)

Convergence rates

Definition (Convergence of a sequence)

The sequence $\mathbf{u}^1, \mathbf{u}^2, ..., \mathbf{u}^k, ...$ converges to \mathbf{u}^* (denoted $\lim_{k \to \infty} \mathbf{u}^k = \mathbf{u}^*$), if

$$\forall \ \varepsilon > 0, \exists \ K \in \mathbb{N} : k \ge K \Rightarrow \|\mathbf{u}^k - \mathbf{u}^\star\| \le \varepsilon$$

Convergence rates: the "speed" at which a sequence converges

sublinear: if there exists c > 0 such that

$$\|\mathbf{u}^k - \mathbf{u}^\star\| = O(k^{-c})$$

▶ **linear:** if there exists $\alpha \in (0,1)$ such that

$$\|\mathbf{u}^k - \mathbf{u}^\star\| = O(\alpha^k)$$

Q-linear: if there exists a constant $r \in (0,1)$ such that

$$\lim_{k \to \infty} \frac{\|\mathbf{u}^{k+1} - \mathbf{u}^{\star}\|}{\|\mathbf{u}^k - \mathbf{u}^{\star}\|} = r$$

- **superlinear:** if r = 0, we say that the sequence converges *superlinearly*.
- ▶ quadratic: if there exists a constant $\mu > 0$ such that $\lim_{k \to \infty} \frac{\|\mathbf{u}^{k+1} \mathbf{u}^{\star}\|}{\|\mathbf{u}^k \mathbf{u}^{\star}\|^2} = \mu$

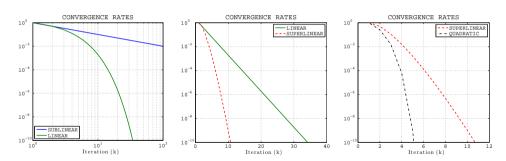
Example: Convergence rates

Examples of sequences that all converge to $u^{\star}=0$:

▶ Sublinear: $u^k = 1/k$

▶ Superlinear: $u^k = k^{-k}$

• Quadratic: $u^k = 0.5^{2^k}$



Wrap up!

- o Next handout will have rate examples!
- o See advanced material for material beyond convexity!
 - ► Star-convexity
 - Invexity
- o Lecture on Monday!

*Jacobian of the self-attention module [5]

Example

We consider the Jacobian of $f: \mathbf{X} \mapsto \sigma_s \left(\mathbf{X} \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{X}^\top \right) \mathbf{X} \mathbf{W}_V^\top$, where σ_s is row-wise softmax, $\mathbf{X} \in \mathbb{R}^{d_s \times d}$, $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_m \times d}$, $f(\mathbf{X}) \in \mathbb{R}^{d_s \times d_m}$.

 $lackbox{lack}$ Define $eta_i := \sigma_s \left(\mathbf{X}^{(i,:)} \mathbf{W}_Q^ op \mathbf{W}_K \mathbf{X}^ op
ight)^ op \in \mathbb{R}^{d_s}$ We can reformulate the definition above as:

$$f(\mathbf{X}) = egin{bmatrix} oldsymbol{eta}_1^{ op} \ dots \ oldsymbol{eta}_{d_c}^{ op} \end{bmatrix} \mathbf{X} \mathbf{W}_V^{ op}.$$

By the product rule:

$$\frac{\partial f(\mathbf{X})}{\partial X^{(p,k)}} = \begin{bmatrix} \frac{\partial \boldsymbol{\beta}_{1}^{\top}}{\partial X^{(p,k)}} \\ \vdots \\ \frac{\partial \boldsymbol{\beta}_{d_{s}}^{\top}}{\partial X^{(p,k)}} \end{bmatrix} \mathbf{X} \mathbf{W}_{V}^{\top} + \begin{bmatrix} \boldsymbol{\beta}_{1}^{\top} \\ \vdots \\ \boldsymbol{\beta}_{d_{s}}^{\top} \end{bmatrix} \frac{\partial (\mathbf{X} \mathbf{W}_{V}^{\top})}{\partial X^{(p,k)}}. \tag{1}$$

*Jacobian of self-attention module [5]

▶ Suppose $\beta = \mathsf{Softmax}(\mathbf{u}) \in \mathbb{R}^{d_s}$, then $\frac{\partial \beta}{\partial \mathbf{u}} = \mathsf{diag}(\beta) - \beta \beta^{\top}$. This is because:

$$\blacktriangleright \text{ We can reformulate } \pmb{\beta} \text{ as: } \pmb{\beta} = \begin{bmatrix} \frac{\exp{(u^{(1)})}}{\sum_{i=1}^{d_s} \exp{(u^{(i)})}} \\ \vdots \\ \frac{\exp{(u^{(d_s)})}}{\sum_{i=1}^{d_s} \exp{(u^{(i)})}} \end{bmatrix}.$$

Thus

$$\frac{\partial \beta^{(j)}}{\partial u^{(k)}} = \frac{\partial \frac{\exp{(u^{(j)})}}{\sum_{i=1}^{d_s} \exp{(u^{(i)})}}}{\partial u^{(k)}} = \begin{cases} \frac{-\exp{(u^{(j)})} - \exp{(u^{(k)})}}{(\sum_{i=1}^{d_s} \exp{(u^{(i)})})^2} & \text{if } j \neq k \\ \frac{\exp{(u^{(k)})} \sum_{i=1}^{d_s} \exp{(u^{(i)})} - (\exp{(u^{(k)})})^2}{(\sum_{i=1}^{d_s} \exp{(u^{(i)})})^2} & \text{if } j = k \end{cases}$$

$$= \begin{cases} -\beta^{(j)} \beta^{(k)} & \text{if } j \neq k \\ \beta^{(k)} - \beta^{(j)} \beta^{(k)} & \text{if } j = k \end{cases}.$$

► Thus

$$\frac{\partial \boldsymbol{\beta}}{\partial \mathbf{n}} = \operatorname{diag}(\boldsymbol{\beta}) - \boldsymbol{\beta} \boldsymbol{\beta}^{\top}. \tag{2}$$

*Jacobian of self-attention module [5]

▶ Then we can calculate the term $\frac{\partial \beta_i}{\partial X^{(p,k)}}$ for $i \in [d_s]$ in the first part of Eq. (1).

$$\frac{\partial \boldsymbol{\beta}_{i}}{\partial X^{(p,k)}} = \left(\operatorname{diag}(\boldsymbol{\beta}_{i}) - \boldsymbol{\beta}_{i} \boldsymbol{\beta}_{i}^{\top}\right) \frac{\partial \left(\mathbf{X} \mathbf{W}_{K}^{\top} \mathbf{W}_{Q} \mathbf{X}^{(i,:)^{\top}}\right)}{\partial X^{(p,k)}} \\
= \left(\operatorname{diag}(\boldsymbol{\beta}_{i}) - \boldsymbol{\beta}_{i} \boldsymbol{\beta}_{i}^{\top}\right) \left(\boldsymbol{e}_{p} \boldsymbol{e}_{k}^{\top} \mathbf{W}_{K}^{\top} \mathbf{W}_{Q} \mathbf{X}^{(i,:)^{\top}} + \mathbf{X} \mathbf{W}_{K}^{\top} \mathbf{W}_{Q} \boldsymbol{e}_{k} \delta_{ip}\right), \tag{3}$$

where e_p is the p^{th} canonical basis vector of \mathbb{R}^{d_s} , e_k is the k^{th} canonical basis vector of \mathbb{R}^d .

Next, let's consider the second term in Eq. (1):

$$\frac{\partial (\mathbf{X} \mathbf{W}_{V}^{\top})}{\partial X^{(p,k)}} = e_{p} e_{k}^{\top} \mathbf{W}_{V}^{\top}. \tag{4}$$

Lastly, substituting Eq. (3) and Eq. (4) into Eq. (1):

$$\frac{\partial f(\mathbf{X})}{\partial X^{(p,k)}} = \begin{bmatrix} \left(\operatorname{diag}(\beta_1) - \beta_1 \beta_1^\top \right) \left(e_p e_k^\top \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{X}^{(1,:)^\top} + \mathbf{X} \mathbf{W}_K^\top \mathbf{W}_Q e_k \delta_{1p} \right) \\ \vdots \\ \left(\operatorname{diag}(\beta_{d_s}) - \beta_{d_s} \beta_{d_s}^\top \right) \left(e_p e_k^\top \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{X}^{(d_s,:)^\top} + \mathbf{X} \mathbf{W}_K^\top \mathbf{W}_Q e_k \delta_{d_sp} \right) \end{bmatrix} \mathbf{X} \mathbf{W}_V^\top + \begin{bmatrix} \beta_1^\top \\ \vdots \\ \beta_{d_s}^\top \end{bmatrix} e_p e_k^\top \mathbf{W}_V^\top.$$

State Space Model (SSM)

- o A state space model represents a system based on a set of first-order differential equations.
- \circ State equation is given by $\dot{\mathbf{x}} = f(\mathbf{x}, \mathbf{u})$, where
 - $\mathbf{x} \in \mathbb{R}^p$ is the state vector,
 - $\mathbf{u} \in \mathbb{R}^n$ is the input vector,
 - ▶ *f* is a (potentially nonlinear) function.
- \circ Output equation is given by $\mathbf{y}=g(\mathbf{x},\mathbf{u})$, where
 - $\mathbf{y} \in \mathbb{R}^d$ is the output vector.
 - q is a function.
- o For linear systems, the state equation and output equation are given by

$$\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$$

$$y = Cx + Du$$

where

- $\mathbf{A} \in \mathbb{R}^{p \times p}$ is the state matrix.
- $ightharpoonup \mathbf{B} \in \mathbb{R}^{p \times n}$ is the input matrix.

- $ightharpoonup \mathbf{C} \in \mathbb{R}^{d imes p}$ is the output matrix.
- $ightharpoonup \mathbf{D} \in \mathbb{R}^{d \times n}$ is the feedthrough matrix.

Jacobian of SSM

Example

Consider a nonlinear system $\ddot{\theta} = -\sin(\theta) + \mathbf{u}$. Define $\mathbf{x}_1 = \dot{\theta}$ and $\mathbf{x}_2 = \ddot{\theta}$. Then it can be written as an SSM:

$$\begin{split} \dot{\mathbf{x}} &= \begin{bmatrix} \dot{\mathbf{x}}_1 \\ \dot{\mathbf{x}}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{x}_2 \\ -\sin(\mathbf{x}_1) + \mathbf{u} \end{bmatrix} \\ \mathbf{y} &= \mathbf{x}_1 \end{split}$$

Jacobian Matrices:

$$\mathbf{A} = \frac{\partial f}{\partial \mathbf{x}} = \begin{bmatrix} 0 & 1\\ -\cos(\mathbf{x}_1) & 0 \end{bmatrix}$$
$$\mathbf{B} = \frac{\partial f}{\partial \mathbf{u}} = \begin{bmatrix} 0\\ 1 \end{bmatrix}$$

$$\mathbf{C} = \frac{\partial g}{\partial \mathbf{x}} = \begin{bmatrix} 1 & 0 \end{bmatrix}$$

$$\mathbf{D} = \frac{\partial g}{\partial \mathbf{u}} = 0$$

▶ The nonlinear system can be approximated by a linear system. Given an equilibrium point $(\mathbf{x}_0, \mathbf{u}_0)$, the linearized state and output equations are given by:

$$\delta \dot{\mathbf{x}} \approx \mathbf{A}|_{(\mathbf{x}_0, \mathbf{u}_0)} \delta \mathbf{x} + \mathbf{B}|_{(\mathbf{x}_0, \mathbf{u}_0)} \delta \mathbf{u}$$

$$\delta \mathbf{y} \approx \left. \mathbf{C} \right|_{(\mathbf{x}_0, \mathbf{u}_0)} \delta \mathbf{x} + \left. \mathbf{D} \right|_{(\mathbf{x}_0, \mathbf{u}_0)} \delta \mathbf{u}$$

where $\delta \mathbf{x}$ and $\delta \mathbf{u}$ represent small changes around $(\mathbf{x}_0, \mathbf{u}_0)$. Note that $\dot{\mathbf{x}} = \delta \dot{\mathbf{x}}$ since \mathbf{x}_0 is a constant.

Convex hull

Definition (Convex hull)

Let $\mathcal{Q} \subseteq \mathbb{R}^p$ be a set. The convex hull of \mathcal{Q} , i.e., $conv(\mathcal{Q})$, is the *smallest* convex set that contains \mathcal{Q} .

Definition (Convex hull of points)

Let $\mathcal{Q} \subseteq \mathbb{R}^p$ be a finite set of points with cardinality $|\mathcal{Q}|$. The convex hull of \mathcal{Q} is the set of all convex combinations of its points, i.e.,

$$\mathrm{conv}(\mathcal{Q}) = \left\{ \sum_{i=1}^{|\mathcal{Q}|} \alpha_i \mathbf{x}_i \ : \ \sum_{i=1}^{|\mathcal{Q}|} \alpha_i = 1, \ \alpha_i \geq 0, \forall i, \ \mathbf{x}_i \in \mathcal{Q} \right\}.$$



Figure: (Left) Discrete set of points Q. (Right) Convex hull conv(Q).

*Star convex sets

Definition

 $\mathcal{Q} \subseteq \mathbb{R}^p$ is a *star-shaped* set if there exists a $\mathbf{x}_1 \in \mathcal{Q}$ such that

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

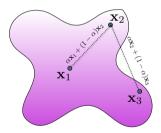


Figure: Example of a star-shaped but not convex set.

*Star convexity

Definition

A function $f \colon \mathcal{Q} \to \mathbb{R} \cup \{+\infty\}$ is called *star-convex* on \mathcal{Q} if there exists a global minimum $\mathbf{x}^\star \in \mathcal{Q}$ such that

$$f(\alpha \mathbf{x}^* + (1 - \alpha)\mathbf{x}) \le \alpha f(\mathbf{x}^*) + (1 - \alpha)f(\mathbf{x}) \quad \forall \mathbf{x} \in \mathcal{Q}, \quad \forall \alpha \in [0, 1].$$

Remarks:

- Any convex function is star-convex.
- \circ Star-convexity can be viewed as convexity between any point ${f x}$ and a global minimum ${f x}^\star.$
- \circ Allows the negative gradient $-\nabla f(\mathbf{x})$ to the desired minimization direction.
- o Consider the following objective function:

$$\min_{\mathbf{x}} f(\mathbf{x}) := \frac{1}{n} \left(\sum_{i=1}^{n} |b_i - \langle \mathbf{a}_i, \mathbf{x} \rangle|^q \right)^{1/q}.$$

- ▶ Star-convex for any real number q when $n \leq p$.
- ▶ Convex for $q \ge 1$.
- ightharpoonup (q=1): the least-absolute deviation estimator. (q=2): the least-squares estimator.

*Invex function

Definition

Let $\mathcal Q$ be an open set in $\mathbb R^p$. A differentiable function $f\colon \mathcal Q\to\mathbb R$ is called *invex* if there exists a function $\eta:\mathcal Q\times\mathcal Q\to\mathbb R^p$ such that

$$f(\mathbf{x}) \ge f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \ \eta(\mathbf{x}, \mathbf{y}) \rangle, \quad \forall \mathbf{x}, \ \mathbf{y} \in \mathcal{Q}.$$

Remarks: o Any convex function is invex function: $\eta(\mathbf{x}, \mathbf{y}) = \mathbf{x} - \mathbf{y}$.

• Any local minima in an invex function is global minima!

Proof: \circ Suppose \mathbf{x}^* is a local minimum, then $\nabla f(\mathbf{x}^*) = 0$. By the definition above, we have

$$f(\mathbf{x}) \ge f(\mathbf{x}^*) + \langle 0, \eta(\mathbf{x}, \mathbf{y}) \rangle = f(\mathbf{x}^*), \quad \forall \mathbf{x} \in \mathcal{Q}.$$

 $\circ \Rightarrow \mathbf{x}^*$ is also a global minimum.

Example (Causality via directed acyclic graph (DAG) learning [1])

For any s>0, define $f^s: \{ \boldsymbol{X} \in \mathbb{R}^{d \times d} \mid s> \rho(\boldsymbol{X} \circ \boldsymbol{X}) \} \to \mathbb{R}$ as $f^s(\boldsymbol{X}) \stackrel{\text{def}}{=} -\log \det(s\boldsymbol{I} - \boldsymbol{X} \circ \boldsymbol{X}) + d\log s$, where \circ is the Hadamard product, $\rho(\cdot)$ is the spectral radius, and \boldsymbol{X} is the graph weighted adjacency matrix.

▶ Then, f^s is an invex function. $f^s(X) \ge 0$ with $f^s(X) = 0$ if and only if X is a DAG.

*Self-concordant functions [4]

Definition (Self-concordant functions in 1-dimension)

A convex function $\varphi:\mathbb{R}\to\mathbb{R}$ is self-concordant if

$$|\varphi'''(t)| \le 2\varphi''(t)^{3/2}, \quad \forall t \in \mathbb{R}.$$

*Self-concordant functions [4]

Definition (Self-concordant functions in 1-dimension)

A convex function $\varphi:\mathbb{R}\to\mathbb{R}$ is self-concordant if

$$|\varphi'''(t)| \le 2\varphi''(t)^{3/2}, \quad \forall t \in \mathbb{R}.$$

Affine Invariance of self-concordant functions

Let $\tilde{\varphi}(t) = \varphi(\alpha t + \beta)$ where $\alpha \neq 0$. Then, $\tilde{\varphi}$ is self-concordant iff φ is.

*Self-concordant functions [4]

Definition (Self-concordant functions in 1-dimension)

A convex function $\varphi:\mathbb{R}\to\mathbb{R}$ is self-concordant if

$$|\varphi'''(t)| \le 2\varphi''(t)^{3/2}, \quad \forall t \in \mathbb{R}.$$

Affine Invariance of self-concordant functions

Let $\tilde{\varphi}(t) = \varphi(\alpha t + \beta)$ where $\alpha \neq 0$. Then, $\tilde{\varphi}$ is self-concordant iff φ is.

Important remarks of self-concordance

- 1. Generalize to higher dimension: A convex function $f:\mathbb{R}^n\to\mathbb{R}$ is said to be (standard) self-concordant if $|\varphi'''(t)|\leq 2\varphi''(t)^{3/2}$, where $\varphi(t):=f(\mathbf{x}+t\mathbf{v})$ for all $t\in\mathbb{R}$, $\mathbf{x}\in\mathrm{dom}\,f$ and $\mathbf{v}\in\mathbb{R}^n$ such that $\mathbf{x}+t\mathbf{v}\in\mathrm{dom}\,f$.
- 2. Affine invariance still holds in high dimension.
- 3. Self-concordant functions are efficiently minimized by the Newton method and its variants.

References |

[1] Kevin Bello, Bryon Aragam, and Pradeep Ravikumar.

Dagma: Learning dags via m-matrices and a log-determinant acyclicity characterization.

In Advances in Neural Information Processing Systems, 2022. (Cited on page 64.)

[2] S. Boyd, S.P. Boyd, L. Vandenberghe, and Cambridge University Press. Convex Optimization.

Berichte über verteilte messysteme. Cambridge University Press, 2004.

(Cited on pages 33, 34, 35, and 36.)

[3] Sham M Kakade and Jason D Lee.

Provably correct automatic sub-differentiation for qualified programs.

In Advances in Neural Information Processing Systems, volume 31, 2018.

(Cited on page 40.)

[4] Yurii Nesterov and Arkadii Nemirovskii. Interior-point polynomial algorithms in convex programming. SIAM, 1994.

(Cited on pages 65, 66, and 67.)

References II

[5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin.

Attention is all you need.

In Advances in Neural Information Processing Systems, 2017.

(Cited on pages 56, 57, and 58.)