

IFT 6085 - Lecture 3

Gradients for smooth and for strongly convex functions

This version of the notes has not yet been thoroughly checked. Please report any bugs to the scribes or instructor.

Scribes

Winter 2020: Thong (Bob) Vo

Winter 2019: Amir Raza, Philippe Lacaille, Jonathan Plante

Winter 2018: Philippe Brouillard, Massimo and Lucas Caccia

Instructor: Ioannis Mitliagkas

1 Summary

In the previous lecture we covered the notions of convexity as well as Lipschitz continuity. After introducing these concepts, a bound on the convergence rate of gradient descent of a convex and L -Lipschitz function was demonstrated to scale with the \sqrt{T} .

Building on some of the previous lecture notions, we will introduce guarantees on the convergence rate of gradient descent for a stronger family of functions (using stronger assumptions), namely β -smooth and α -strong convex functions.

2 Gradient Descent for smooth functions

Definition 1 (β -smoothness). We say that a continuously differentiable function f is β -smooth if its gradient ∇f is β -Lipschitz, that is

$$\|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\| \quad \Rightarrow \quad \|f(x) - f(y)\| \leq \frac{1}{2} \beta \|x - y\|^2$$

If we recall Lipschitz continuity from Lecture 2, simply speaking, an L -Lipschitz function is limited by how *quickly* its output can change. By imposing this same constraint on the gradients of a function, β -smoothness implies they cannot change abruptly and must be bounded by some value as defined above.

In other other words, β -smoothness is putting an upper bound on the curvature of the function. This is equivalent to the eigenvalues of the Hessian being less than β . Note that there can be β -smooth functions which are not twice differentiable. One key benefit of these functions is that their gradients tend to decay when x gets closer to the minimum. In contrast, non-smooth functions may have abrupt bends at the minimum, which cause significant oscillations for gradient descent. Figure 1 illustrates this point by comparing the two scenarios.

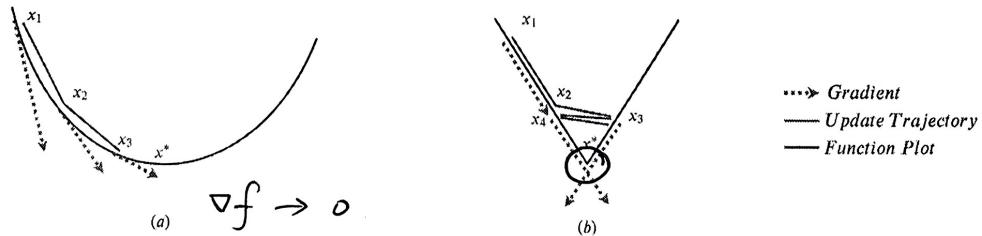


Figure 1: (a) A smooth function with decaying updates. (b) A non-smooth function with oscillating updates.

2.1 Convex and smooth functions

Here we introduce a bound on the convergence rate of a convex and β -smooth function.

Lemma 2 (Quadratic bounds). *Let f be β -smooth on \mathbb{R}^n . Then for any $x, y \in \mathbb{R}^n$, one has*

$$|f(x) - f(y) - \nabla f(y)^T(x - y)| \leq \frac{\beta}{2} \|x - y\|^2 \quad (1)$$

Proof. Let $g : \mathbb{R} \rightarrow \mathbb{R}$ by the rule $g(t) \triangleq f(y + t(x - y))$. The following holds true:

$$f(x) = g(1) \text{ and } f(y) = g(0) \quad f(x) \leq f(y) + \nabla f(y)^T(x - y)$$

Then, we also observe that:

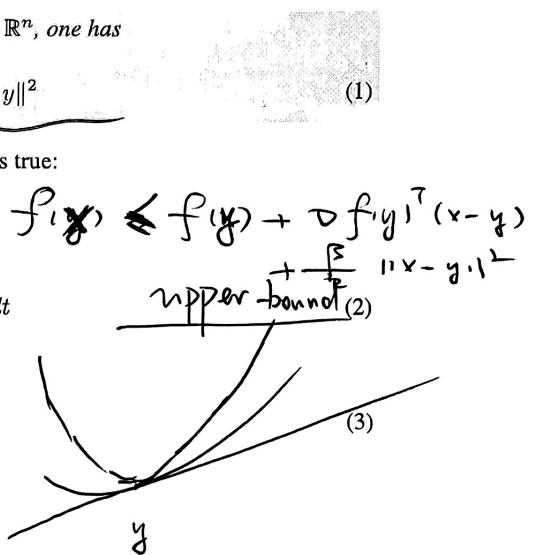
$$f(x) - f(y) = g(1) - g(0) = \int_0^1 g'(t) dt$$

Taking derivative of $g(t)$:

$$g'(t) = \nabla f(y + t(x - y))^T(x - y)$$

Plugging in Equation 2, and 3 to the LHS of 1:

$$\begin{aligned} & |f(x) - f(y) - \nabla f(y)^T(x - y)| \\ &= \left| \int_0^1 \nabla f(y + t(x - y))^T(x - y) dt - \nabla f(y)^T(x - y) \right| \\ &\leq \int_0^1 |\nabla f(y + t(x - y)) - \nabla f(y)|^T(x - y) dt \\ &\leq \int_0^1 \|\nabla f(y + t(x - y)) - \nabla f(y)\| \cdot \|x - y\| dt \quad (\text{applying Cauchy-Schwarz inequality}) \\ &\leq \int_0^1 \beta t \|x - y\|^2 dt \quad (\text{the gradient } \nabla f \text{ is } \beta\text{-Lipschitz.}) \\ &= \frac{\beta}{2} \|x - y\|^2 \end{aligned}$$



□

Lemma 3. *Let f be such that $0 \leq f(x) - f(y) - \nabla f(y)^T(x - y) \leq \frac{\beta}{2} \|x - y\|^2$. Then for any $x, y \in \mathbb{R}^n$, one has*

$$f(x) - f(y) \leq \nabla f(x)^T(x - y) - \frac{1}{2\beta} \|\nabla f(x) - \nabla f(y)\|^2$$

Proof. See [1] Lemma 3.5 for proof. □

Theorem 4. *Let f be convex and β -smooth on \mathbb{R}^n . Then the gradient descent with step size $\gamma = 1/\beta$ satisfies*

$$f(x_k) - f(x^*) \leq \frac{2\beta \|x_1 - x^*\|^2}{k-1} \sim O(\frac{1}{k}) \quad \text{Convergence}$$

Proof. See [1] page 268. □ Vore

In comparison to the convergence analysis for a convex and L -Lipschitz function, the following points of improvement are observed in Theorem 4 over the previous result:

- No averaging of terms

- x_k is a good solution, i.e. no reference to previous iterations
- Convergence scales linearly with number of steps, convergence rate of order $O(1/T)$ compared to $O(1/\sqrt{T})$
- Ideal step size γ is constant and does not depend on T (number of steps taken).

As an observation of $\|x_1 - x^*\|^2$, the bound will be tighter if x_1 will be closer to x^* (minima), and looser if x_1 will be farther. Bounds that don't depend on the initial value of x will be discussed later in this lecture.

3 Strong convexity

Definition 5 (Strong convexity). A function $f(x)$ is α -strongly convex, if for $\alpha > 0$, $\forall x \in \text{dom}(f)$,

$$f(x) - \frac{\alpha}{2} \|x\|^2 \text{ is convex.}$$

Strong convexity provides a lower bound for the function's curvature. The function must have strictly positive curvature. In other words, all eigenvalues of the Hessian of a α -strongly convex function are lower bounded by α . We can write this in terms of positive-semi definiteness as

$$\nabla^2 f(x) \succeq \alpha I \iff \nabla^2 f(x) - \alpha I \succeq 0$$

For example, $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = \frac{h}{2}x^2$ is h -strongly convex, but not $(h + \epsilon)$ -strongly convex for $\epsilon > 0$. Figure 2 illustrates examples of two convex functions, of which only one is strongly convex.

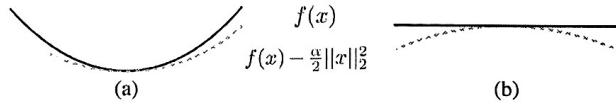


Figure 2: (a) A convex function which is also strongly convex. (b) A convex function which is not strongly convex.

3.1 Strongly convex and Lipschitz functions

Theorem 6. Let f be α -strongly convex and L -Lipschitz. Then the projected subgradient descent after T steps with $\gamma_k = \frac{2}{\alpha(k+1)}$ satisfies

$$f\left(\sum_{k=1}^T \frac{2k}{T(T+1)} x_k\right) - f(x^*) \leq \frac{2L^2}{\alpha(T+1)} \sim O\left(\frac{1}{T}\right)$$

Proof. See [1] page 277 □

With Theorem 6, we can notice how moving from convexity to strong convexity may affect the convergence rate of gradient descent. We previously tackled convexity paired with Lipschitz continuity in Lecture 2 and the similarity with the new convergence guarantees is pretty noticeable. By moving to strong-convexity, we can notice a few things:

- Still have averaging of terms, where $\sum_{k=1}^T \frac{2k}{T(T+1)} x_k$ is a non-uniform averaging scheme with more recent iterates having more weight $\rightarrow \frac{1}{T} \sum_{k=1}^T \left(\frac{2k}{T+1}\right) x_k$ weights.
- Current x_k solution is not appropriate, i.e. ideally only evaluate at x_k
- Convergence now scales linearly with number of steps
- Although not constant, the step size γ_k is diminishing at every step

$$\gamma_k \rightarrow 0. \quad 3$$

Unlike previous results, the right hand side does not have any terms dependent on distance from x^* here. Intuitively this is a consequence of two conditions which have to be met for a Strong Convex function:

- Norm of Gradient increases as we go further away from minima.
- Gradient for a L-Lipschitz function is bounded, and can not keep increasing.

Thus to have a finite bound, no distance term on the right hand side of the bound.

One could hope that by combining strong convexity along with smoothness, gradient descent may present stronger convergence guarantees. We will see how smoothness removes dependency from the averaging scheme.

Note: α Strongly convex and L-Lipschitz condition is a special case because the upper bound L-Lipschitz condition will ultimately conflict with the lower bound α Strongly convex grow rate. Therefore, such functions are typically defined in a range, e.g. $x \in [-1, 1]$.

3.2 Strongly convex and smooth functions

Recalling Lemma 2 (Quadratic bounds), it tells us that a β -smooth function f is sandwiched between 2 quadratics because of the following inequality:

$$f(y) + \nabla f(y)^T(x - y) - \frac{\beta}{2}\|x - y\|^2 \leq f(x) \leq f(y) + \nabla f(y)^T(x - y) + \frac{\beta}{2}\|x - y\|^2 \quad (4)$$

Now we introduce another lemma which allows us to lower bound an α -strong convex function.

Lemma 7. Let f be λ -strongly convex. Then $\forall x, y$, we have:

$$f(y) - f(x) \leq \nabla f(y)^T(y - x) - \frac{\lambda}{2}\|x - y\|^2$$

The strong convexity parameter λ is a measure of the curvature of f .

By rearranging terms, this tells us that a λ -strong convex function can be lower bounded by the following inequality:

$$f(x) \geq f(y) - \nabla f(y)^T(y - x) + \frac{\lambda}{2}\|x - y\|^2 \quad (5)$$

Figure 3 showcases the resulting bounds from both the smoothness and the strong convexity constraints. The shaded area in each sub-figure is showing the area validated by the respective bound(s).

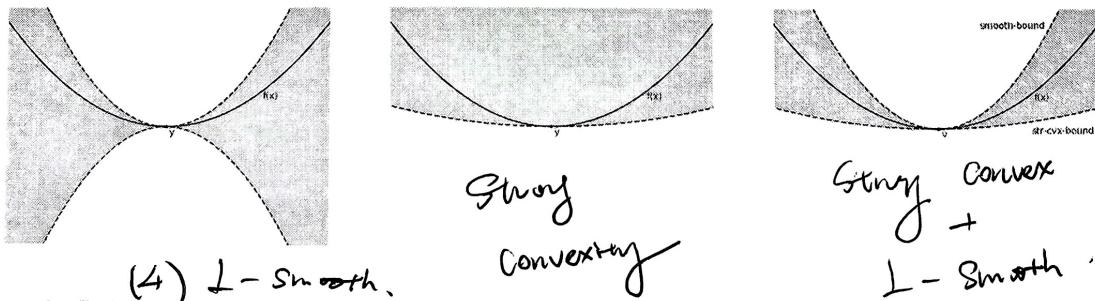


Figure 3: (Left) Upper and lower bounds from equation 4 for smoothness constraint (Middle) Lower bound from equation 5 for strong convexity constraint (Right) Combination of upper bound from smoothness and lower bound from strong convexity

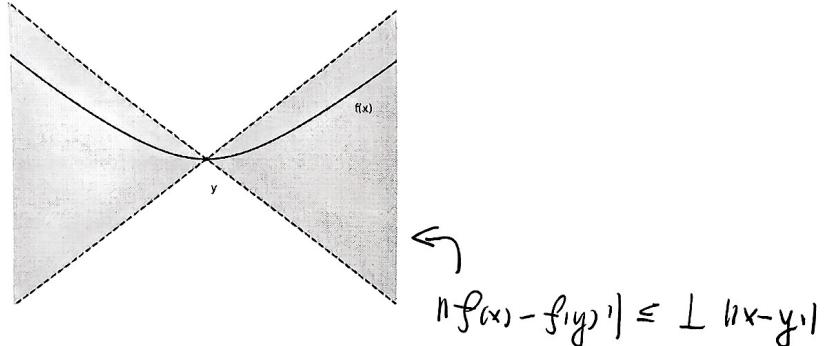


Figure 4: For an L_f -Lipschitz continuous function, the green region shows where the function would exist. We can imagine that without smoothness and only L -Lipschitz in equation 4, the accepted region would be having linear boundaries

Lemma 8 (Coercivity of the gradient). *Let f be β -smooth and λ -strongly convex. Then for all x and y , we have:*

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{\lambda\beta}{\lambda + \beta} \|x - y\|^2 + \frac{1}{\lambda + \beta} \|\nabla f(x) - \nabla f(y)\|^2$$

Proof. See [1] Lemma 3.11 for proof. □

Theorem 9. For f a λ -strongly convex and β -smooth function, gradient descent with $\gamma = \frac{2}{\lambda + \beta}$ satisfies:

$$f(x_{k+1}) - f(x^*) \leq \frac{\beta}{2} \exp\left(-\frac{4k}{\kappa + 1}\right) \|x_1 - x^*\|^2$$

where $\kappa \triangleq \frac{\beta}{\lambda}$ is the condition number. Passes $\frac{1}{\frac{1}{\kappa + 1}} = \boxed{\frac{1}{\kappa + 1}}$ $\boxed{\frac{e^{-\frac{1}{\kappa + 1}}}{\frac{1}{\kappa + 1}}} = \boxed{x e^{-\frac{1}{\kappa + 1}}} \quad \boxed{x \rightarrow 0.}$

Proof. (Theorem 9) First, let's define the distance between the coordinate at the iteration k and the optimal point as: \approx

$$D_k \triangleq \|x_k - x^*\|$$

Then,

$$D_{k+1}^2 = \|x_{k+1} - x^*\|^2$$

By replacing x_{k+1} by its definition $x_k - \gamma \nabla f(x_k)$, we get:

$$D_{k+1}^2 = \|x_k - \gamma \nabla f(x_k) - x^*\|^2$$

By expanding the square, we get:

$$D_{k+1}^2 = \|x_k - x^*\|^2 - 2\gamma \langle \nabla f(x_k), x_k - x^* \rangle + \gamma^2 \|\nabla f(x_k)\|^2$$

By doing a slight modification to the second term, we can apply the lemma of coercivity of the gradient. Since $\nabla f(x^*) = 0$ this equality holds:

$$\langle \nabla f(x_k), x_k - x^* \rangle = \langle \nabla f(x_k) - \nabla f(x^*), x_k - x^* \rangle$$

And by applying the lemma of coercivity of the gradient, we get an upper bound: Strongly convex.

$$\langle \nabla f(x_k) - \nabla f(x^*), x_k - x^* \rangle \geq \frac{\lambda\beta}{\lambda + \beta} D_k^2 + \frac{1}{\lambda + \beta} \|\nabla f(x_k) - \nabla f(x^*)\|^2$$

By replacing the second term by the upper bound we just found, we get:

$$D_{k+1}^2 \leq D_k^2 - 2\gamma \left(\frac{\lambda\beta}{\lambda + \beta} D_k^2 + \frac{1}{\lambda + \beta} \|\nabla f(x_k)\|^2 \right) + \gamma^2 \|\nabla f(x_k)\|^2$$

We can rearrange the terms and add $-\nabla f(x^*)$ again inside the norm terms:

$$D_{k+1}^2 \leq \left(1 - \frac{2\gamma\lambda\beta}{\lambda + \beta}\right) D_k^2 + \left(\frac{-2\gamma}{\lambda + \beta} + \gamma^2\right) \|\nabla f(x_k) - \nabla f(x^*)\|^2$$

The term $\left(1 - \frac{2\gamma\lambda\beta}{\lambda + \beta}\right)$ is useful, because we can show that it is less than 1 and has geometric convergence. Further steps will try to simplify the remaining terms.

We can change the norm term by using the fact that f is β -smooth:

$$D_{k+1}^2 \leq \left(1 - \frac{2\gamma\lambda\beta}{\lambda + \beta}\right) D_k^2 + \left(\frac{-2\gamma}{\lambda + \beta} + \gamma^2\right) \beta^2 D_k^2$$

Let $\gamma = \frac{2}{\lambda + \beta}$, then:

$$D_{k+1}^2 \leq \left(1 - \frac{4\lambda\beta}{(\lambda + \beta)^2}\right) D_k^2$$

By unrolling the recursion and since $\left(\frac{\kappa-1}{\kappa+1}\right)^2 = \left(1 - \frac{4\lambda\beta}{(\lambda + \beta)^2}\right)$ we get:

$$D_{k+1}^2 \leq \left(\frac{\kappa-1}{\kappa+1}\right)^{2k} D_1^2$$

Since $\exp(-x) \geq 1 - x$ for every x , we get:

$$\begin{aligned} D_{k+1}^2 &\leq \exp\left(-\frac{4k}{\kappa+1}\right) \overbrace{D_1^2} \\ \text{By } \beta\text{-smoothness we finally have: } \|x_{k+1} - x^*\|^2 &\leq \downarrow \text{--- Smooth} \\ f(x_{k+1}) - f(x^*) &\leq \frac{\beta}{2} \exp\left(-\frac{4k}{\kappa+1}\right) \|x_1 - x^*\|^2 \end{aligned}$$

□

Theorem 9 observations

- x_k solution is a good solution, i.e. no reference to previous iterations
- Convergence rate of order $O(\exp(-T))$
- κ measures how far apart the upper and lower bounds are (see Figure 3). It can be interpreted as the ratio of largest to smallest curvature of the function.
- The smaller the condition number κ is, the less iterations are required to converge. Intuitively, the accepted region between the bounds will be smaller.
- Consequently, the greater β is the more iterations will be required to converge. This is logical since a constant step size on a function with a steep gradient will cause a greater change in the function value.

4 Comparison of optimization properties for different function classes

The table below summarizes various convergence properties of discussed functions classes. From left to right, the assumptions on properties of these classes increase, from Convex and L-Lipschitz to λ strongly convex and β smooth. In the same direction, the convergence rates are also increase, aligned with stronger assumptions on those functions.

Table 1: Comparison of different function classes

Function Class	cvx, L-Lipschitz	cvx, β smooth	α str-cvx, L-Lipschitz	λ str-cvx, β smooth
Optimal Step size	$\gamma = \frac{\ x_1 - x^*\ _2}{L\sqrt{T}}$	$\gamma = \frac{1}{\beta}$	$\gamma = \frac{2}{\alpha(k+1)}$	$\gamma = \frac{2}{\lambda+\beta}$
Convergence Rate	$\mathcal{O}(1/\sqrt{T})$	$\mathcal{O}(1/T)$	$\mathcal{O}(1/T)$	$\mathcal{O}(\exp(-T))$
Sub-optimal gap	$f\left(\frac{1}{T} \sum_{k=1}^T x_k\right) - f(x^*)$	$(x_k) - f(x^*)$	$f\left(\sum_{k=1}^T \frac{2k}{T(T+1)} x_k\right) - f(x^*)$	$f(x_{k+1}) - f(x^*)$
Bounds of the gap	$\leq \frac{\ x_1 - x^*\ _2^2}{\sqrt{T}}$	$\leq \frac{2\beta\ x_1 - x^*\ ^2}{k-1}$	$\leq \frac{2L^2}{\alpha(T+1)}$	$\leq \frac{\beta}{2} \exp\left(-\frac{4k}{\kappa+1}\right) D_1^2$

References

- [1] S. Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.

数学分析 II

第1次参考解答

2025年3月12日

一、基础题

Critical point

1. 判断 $x = 0$ 是否为以下函数的驻点，是否为以下函数的极值点：

$$(1) f(x) = x^3 - x^2$$

$$(2) f(x) = \begin{cases} x^2 \sin \frac{1}{x}, & (x \neq 0) \\ 0, & (x = 0) \end{cases}$$

解答：

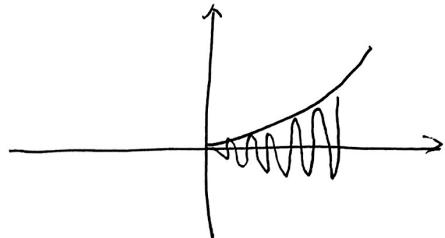
(1) $x = 0$ 是 $f(x)$ 的驻点，也是极大值点

(2) 易知函数 $f(x)$ 在 \mathbb{R} 可导，根据导数定义

$$f'(0) = \lim_{x \rightarrow 0} \frac{f(x) - f(0)}{x - 0} = \lim_{x \rightarrow 0} x \sin \frac{1}{x} = 0$$

所以 $x = 0$ 是 $f(x)$ 的驻点。

但由于可以取到两个趋于0的数列 $a_n = \frac{1}{2n\pi + \frac{\pi}{2}}$, $b_n = \frac{1}{2n\pi + \frac{3\pi}{2}}$, 使得 $f(a_n) > 0, f(b_n) < 0$, 故 $x = 0$ 不是 $f(x)$ 的极值点



2. 判断以下函数在其定义域内的凸性，写出所有拐点。

$$(1) f(x) = x^3 - 4x^2 + x + 6, x \in \mathbb{R}$$

$$(2) f(x) = x - \sin x, x \in (0, 2\pi)$$

$$(3) f(x) = e^x - e^{-x}, x \in \mathbb{R}$$

解答：

$$(1) f''(x) = 6x - 8, \text{ 当 } x \leq \frac{4}{3} \text{ 时 } f''(x) \leq 0, \text{ 当 } x > \frac{4}{3} \text{ 时 } f''(x) > 0$$

故 $f(x)$ 在 $(-\infty, \frac{4}{3})$ 上凸，在 $(\frac{4}{3}, +\infty)$ 下凸，拐点为 $\frac{4}{3}$

$$(2) f''(x) = \sin x, \text{ 当 } x \in (0, \pi] \text{ 时 } f''(x) \geq 0, \text{ 当 } x \in [\pi, 2\pi) \text{ 时 } f''(x) \leq 0$$

故 $f(x)$ 在 $(0, \pi]$ 下凸，在 $[\pi, 2\pi)$ 上凸，拐点为 π

$$(3) f''(x) = e^x - e^{-x}, \text{ 当 } x \in (-\infty, 0] \text{ 时 } f''(x) \leq 0, \text{ 当 } x \in [0, +\infty) \text{ 时 } f''(x) \geq 0$$

故 $f(x)$ 在 $(-\infty, 0]$ 上凸，在 $[0, +\infty)$ 下凸，拐点为 0



3. 已知 $f(x), g(x)$ 为区间 I 上的下凸函数，证明：

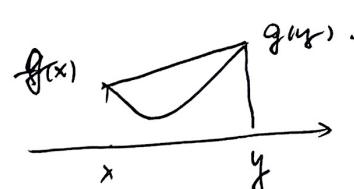
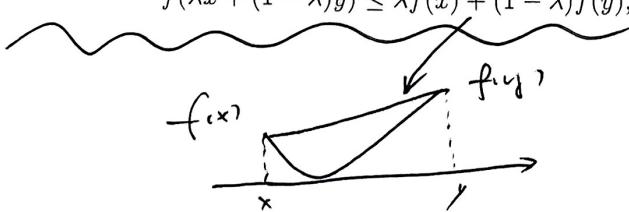
(1) $h(x) = \max\{f(x), g(x)\}$ 为 I 上的下凸函数

(2) 若 $f(x)$ 为严格下凸函数，则 $\phi(x) = e^{f(x)}$ 为 I 上的严格下凸函数

解答：

(1) 根据条件， $\forall x, y \in I, \lambda \in (0, 1)$, 都有

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y), \quad g(\lambda x + (1 - \lambda)y) \leq \lambda g(x) + (1 - \lambda)g(y)$$



而 $h(\lambda x + (1 - \lambda)y) = \begin{cases} f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \leq \lambda h(x) + (1 - \lambda)h(y), & (f_\lambda \geq g_\lambda) \\ g(\lambda x + (1 - \lambda)y) \leq \lambda g(x) + (1 - \lambda)g(y) \leq \lambda h(x) + (1 - \lambda)h(y), & (f_\lambda \leq g_\lambda) \end{cases}$
 f_λ, g_λ 分别表示 $f(\lambda x + (1 - \lambda)y)$ 和 $g(\lambda x + (1 - \lambda)y)$ 。

所以 $\forall x, y \in I, \lambda \in (0, 1)$, 都有 $h(\lambda x + (1 - \lambda)y) \leq \lambda h(x) + (1 - \lambda)h(y)$, 证毕。

(2) 我们证明更一般的结论: 若 $f(x), g(x)$ 是区间 I 上的严格下凸函数, 且 $g(x)$ 在 I 上单增, 则 $\phi(x) = g(f(x))$ 是 I 上的严格下凸函数。证明如下:

已知, $\forall x, y \in I, \lambda \in (0, 1)$, 都有

$$\phi(x) = g(f(x)) \quad \text{Def.} \quad \begin{array}{c} \downarrow \\ \text{Convex} \end{array} \quad \geq 0$$

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y), \quad g(\lambda x + (1 - \lambda)y) < \lambda g(x) + (1 - \lambda)g(y)$$

又因为 $g(x)$ 单增, 则

$$\phi(\lambda x + (1 - \lambda)y) = g(f(\lambda x + (1 - \lambda)y)) \leq g(\lambda f(x) + (1 - \lambda)f(y)) < \lambda g(f(x)) + (1 - \lambda)g(f(y))$$

$\underbrace{\hspace{10em}}$ Def.

所以 $\forall x, y \in I, \lambda \in (0, 1)$

$$\phi(\lambda x + (1 - \lambda)y) < \lambda \phi(x) + (1 - \lambda)\phi(y)$$

证毕。

1. g 个. f convex

2. $g \cdot f$ convex + $f \cdot g$ differentiable.

4. 证明: $f(x)$ 是区间 I 上的下凸函数, 当且仅当 $\forall x_1, x_2 \in I$, $g(\lambda) = f(\lambda x_1 + (1 - \lambda)x_2)$ 是 $[0, 1]$ 上的下凸函数

解答:

先证必要性: 设 $f(x)$ 是区间 I 上的下凸函数, 那么 $\forall \lambda_1, \lambda_2 \in I$ 和 $k \in [0, 1]$, 总有

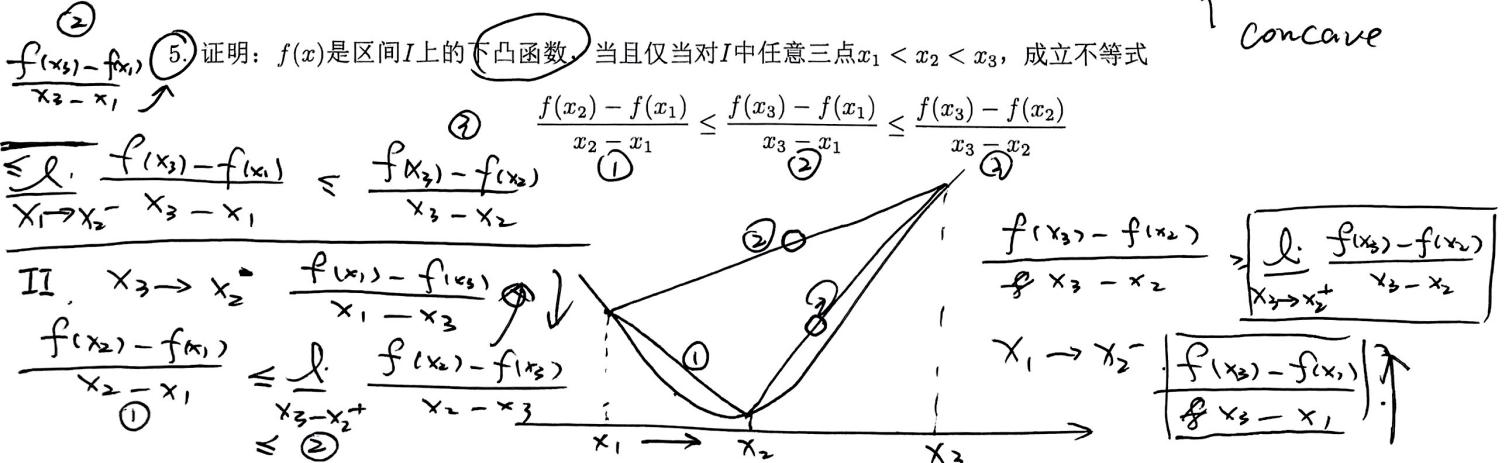
$$\begin{aligned} g(k\lambda_1 + (1 - k)\lambda_2) &= f(\underbrace{k\lambda_1 + (1 - k)\lambda_2}_{\text{Def.}} x_1 + (1 - \underbrace{k\lambda_1 + (1 - k)\lambda_2}_{\text{Def.}})x_2) \\ &= f(k\lambda_1 x_1 + (1 - k)\lambda_2 x_1 + x_2 - k\lambda_1 x_2 - (1 - k)\lambda_2 x_2) \\ &= f(k(\underbrace{\lambda_1 x_1 + (1 - \lambda_1)x_2}_{\text{Def.}}) + (1 - k)(\underbrace{\lambda_2 x_1 + (1 - \lambda_2)x_2}_{\text{Def.}})) \\ &\leq kf(\lambda_1 x_1 + (1 - \lambda_1)x_2) + (1 - k)f(\lambda_2 x_1 + (1 - \lambda_2)x_2)) \\ &= kg(\lambda_1) + (1 - k)g(\lambda_2) \end{aligned}$$

再证充分性: 设 $g(\lambda) = f(\lambda x_1 + (1 - \lambda)x_2)$ 是 $[0, 1]$ 上的下凸函数, 那么 $\forall x_1, x_2 \in I$ 和 $\lambda \in (0, 1)$

$$\begin{aligned} f(\lambda x_1 + (1 - \lambda)x_2) &= g(\lambda) = g(1 + \underbrace{(1 - \lambda)}_{\text{Def.}} \cdot 0) \\ &\leq \lambda g(1) + (1 - \lambda)g(0) = \lambda f(x_1) + (1 - \lambda)f(x_2) \end{aligned}$$

证毕。

Convex
Concave



需要证明的结论可以看作三个不等式：

$$\frac{f(x_2) - f(x_1)}{x_2 - x_1} - \frac{f(x_3) - f(x_1)}{x_3 - x_1} \leq 0$$

$$\frac{f(x_2) - f(x_1)}{x_2 - x_1} - \frac{f(x_3) - f(x_2)}{x_3 - x_2} \leq 0$$

$$\frac{f(x_3) - f(x_1)}{x_3 - x_1} - \frac{f(x_3) - f(x_2)}{x_3 - x_2} \leq 0$$

由于 $x_1 < x_2 < x_3$, 注意到：

$$\begin{aligned} & \begin{vmatrix} 1 & 1 & 1 \\ x_1 & x_2 & x_3 \\ f(x_1) & f(x_2) & f(x_3) \end{vmatrix} = \begin{vmatrix} x_2 - x_1 & x_3 - x_1 \\ f(x_2) - f(x_1) & f(x_3) - f(x_1) \end{vmatrix} \\ & = \begin{vmatrix} x_2 - x_1 & x_3 - x_2 \\ f(x_2) - f(x_1) & f(x_3) - f(x_2) \end{vmatrix} = \begin{vmatrix} x_3 - x_1 & x_3 - x_2 \\ f(x_3) - f(x_1) & f(x_3) - f(x_2) \end{vmatrix} \end{aligned}$$

因此所需证明的三个不等式等价，故只需要证明下凸函数和其中一个不等式等价即可。证明如下：

由于 $x_1 < x_2 < x_3$, 可令

$$\lambda = \frac{x_3 - x_2}{x_3 - x_1}$$

则 λ 满足 $0 < \lambda < 1$ 且 $x_2 = \lambda x_1 + (1 - \lambda)x_3$

先证必要性：由于 f 下凸，有

$$f(x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_3)$$

代入 λ 后稍加整理，得：

$$(x_3 - x_1)[f(x_2) - f(x_1)] \leq (x_2 - x_1)[f(x_3) - f(x_1)]$$

即

$$\frac{f(x_2) - f(x_1)}{x_2 - x_1} \leq \frac{f(x_3) - f(x_1)}{x_3 - x_1}$$

再证充分性：若对 I 中任意三点 $x_1 < x_2 < x_3$

$$\frac{f(x_2) - f(x_1)}{x_2 - x_1} \leq \frac{f(x_3) - f(x_1)}{x_3 - x_1}$$

则

$$f(x_2) \leq \frac{x_3 - x_2}{x_3 - x_1} f(x_1) + \left(1 - \frac{x_3 - x_2}{x_3 - x_1}\right) f(x_3)$$

即

$$f\left(\frac{x_3 - x_2}{x_3 - x_1} x_1 + \left(1 - \frac{x_3 - x_2}{x_3 - x_1}\right) x_3\right) \leq \frac{x_3 - x_2}{x_3 - x_1} f(x_1) + \left(1 - \frac{x_3 - x_2}{x_3 - x_1}\right) f(x_3)$$

故 $\forall x_1, x_3 \in I, \lambda \in (0, 1)$

$$f(\lambda x_1 + (1 - \lambda)x_3) \leq \lambda f(x_1) + (1 - \lambda)f(x_3)$$

证毕。

问题

1. 对于下凸函数，有如下三种定义：称 $f(x)$ 为 I 上的下凸函数，如果

定义 1 (通常定义) $\forall x, y \in I, \lambda \in (0, 1)$, 有 $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$

定义 2 (中点定义) $\forall x_1, x_2 \in I$, 有 $f\left(\frac{x_1 + x_2}{2}\right) \leq \frac{f(x_1) + f(x_2)}{2}$ $\lambda = \frac{1}{2}$

定义 3 (均值点定义) $\forall x_1, x_2, \dots, x_n \in I$, 有 $f\left(\frac{x_1 + x_2 + \dots + x_n}{n}\right) \leq \frac{f(x_1) + f(x_2) + \dots + f(x_n)}{n}$

(1) 证明定义 2 和 定义 3 等价

(2) 当 $f(x)$ 为连续函数时，证明以上三个定义相互等价

(3) 除了以上三种定义，思考其他的下凸函数的等价刻画

解答：

(1) 定义 3 \rightarrow 定义 2 是显然的，只证 $\text{定义 } 2 \rightarrow \text{定义 } 3$ 。这里采用反向归纳法：先证明命题对于所有 $n = 2^k$ 均成立，再证明当 $n = k + 1$ 时命题也成立。

由于当 $n = 2$ 时， $\forall x_1, x_2 \in I$, 有 $f\left(\frac{x_1 + x_2}{2}\right) \leq \frac{f(x_1) + f(x_2)}{2}$ ，那么 $\forall x_1, x_2, x_3, x_4 \in I$, 有

$$\boxed{f\left(\frac{x_1 + x_2 + x_3 + x_4}{4}\right) = f\left(\frac{\frac{x_1 + x_2}{2} + \frac{x_3 + x_4}{2}}{2}\right) \leq \frac{f\left(\frac{x_1 + x_2}{2}\right) + f\left(\frac{x_3 + x_4}{2}\right)}{2} \leq \frac{f(x_1) + f(x_2) + f(x_3) + f(x_4)}{4}}$$

事实上，对于任意的 $n = 2^k$ ，重复以上方法可知

$$f\left(\frac{x_1 + x_2 + \dots + x_{2^k}}{2^k}\right) \leq \frac{f(x_1) + f(x_2) + \dots + f(x_{2^k})}{2^k}$$

于是命题对于所有 $n = 2^k$ 均成立。

若命题对于 $n = k + 1$ 时成立，往证命题对于 $n = k$ 时成立。记 $A = \frac{x_1 + x_2 + \dots + x_k}{k}$ ，那么

$$A = \frac{x_1 + x_2 + \dots + x_k + A}{k+1}$$

由于

$$f(A) = f\left(\frac{x_1 + x_2 + \dots + x_k + A}{k+1}\right) \leq \frac{f(x_1) + f(x_2) + \dots + f(x_k) + f(A)}{k+1}$$

于是

$$(k+1)f(A) \leq f(x_1) + f(x_2) + \dots + f(x_k) + f(A)$$

即

$$kf(A) \leq f(x_1) + f(x_2) + \dots + f(x_k)$$

代入 A 的表达式得

$$\boxed{f\left(\frac{x_1 + x_2 + \dots + x_k}{k}\right) \leq \frac{f(x_1) + f(x_2) + \dots + f(x_k)}{k}} \quad \begin{array}{l} \text{Jensen} \\ \text{inequality} \end{array}$$

所以 $n = k$ 时命题成立，归纳法证毕。

(2) (a) 先证 $\text{定义 } 1 \rightarrow \text{定义 } 2, 3$:

在 $\text{定义 } 1$ 中令 $\lambda = \frac{1}{2}$ ，则

$$f\left(\frac{x+y}{2}\right) \leq \frac{f(x) + f(y)}{2}$$

即为定义2，故定义1本身蕴含定义2。

(b) 再证定义2、3 \rightarrow 定义1：

任取 $x_1, x_2 \in I$ ，先证明当 λ 为有理数时结论成立。不妨设 $\lambda = \frac{m}{n} \in (0, 1)$ ($m < n$ 为自然数)，则

$$\begin{aligned} f(\lambda x_1 + (1 - \lambda)x_2) &= f\left(\frac{m}{n}x_1 + (1 - \frac{m}{n})x_2\right) = f\left(\frac{mx_1 + (n - m)x_2}{n}\right) \\ &\leq \frac{mf(x_1) + (n - m)f(x_2)}{n} = \lambda f(x_1) + (1 - \lambda)f(x_2) \end{aligned}$$

于是 λ 为有理数的情况获证。

若 $\lambda \in (0, 1)$ 为无理数，则 $\exists \lambda_n \in (0, 1)$ ($n = 1, 2, \dots$) 使得 $n \rightarrow +\infty$ 时 $\lambda_n \rightarrow \lambda$ ，从而由 $f(x)$ 的连续性：

$$f(\lambda x_1 + (1 - \lambda)x_2) = f\left(\lim_{n \rightarrow +\infty} (\lambda_n x_1 + (1 - \lambda_n)x_2)\right) = \lim_{n \rightarrow +\infty} f(\lambda_n x_1 + (1 - \lambda_n)x_2)$$

对于有理数 λ_n ，已经证明有

$$f(\lambda_n x_1 + (1 - \lambda_n)x_2) \leq \lambda_n f(x_1) + (1 - \lambda_n)f(x_2)$$

对上式取极限，得

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

这就证明了定义2、3蕴含定义1。

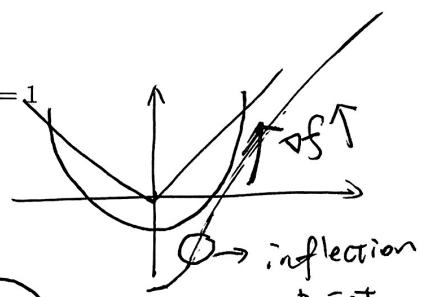
注：从上述证明可以看出，定义1 \rightarrow 定义2、3无须连续性，而定义2、3 \rightarrow 定义1需要连续性，所以定义1强于定义2、3，因此作为我们常用的下凸函数定义。

(3) 请各位小导师们带同学们总结大课讲过一些等价刻画，包括但不限于：

(a) 上方图形 $\{(x, y) \in \mathbb{R} : f(x) \leq y, x \in I\}$ 是凸集

(b) $\forall a \in I, F(x) = \frac{f(x) - f(a)}{x - a}$ ($x \neq a$) 单增

(c) Jensen不等式： $f\left(\sum_{i=1}^n q_i x_i\right) \leq \sum_{i=1}^n q_i f(x_i)$, $\forall n, \forall x_i \in I, \forall q_i \geq 0$ 满足 $\sum_{i=1}^n q_i = 1$



2. $f(x)$ 是定义在 \mathbb{R} 上的函数， $y = |x|$ 是 $f(x)$ 的渐近线，且存在 x_0 使得 $f(x_0) < 0$

(1) 证明： $f(x)$ 不可能是 \mathbb{R} 上的下凸函数

(2) 若 $f(x)$ 在 \mathbb{R} 上可导，且题设“存在 x_0 使得 $f(x_0) < 0$ ”改为“存在 x_0 使得 $f'(x_0) = 0$ ”，证明相同的结论

解答：

(1) 由对称性不妨设 $x_0 \geq 0$ ，因为 $y = |x|$ 是 $f(x)$ 的渐近线，知

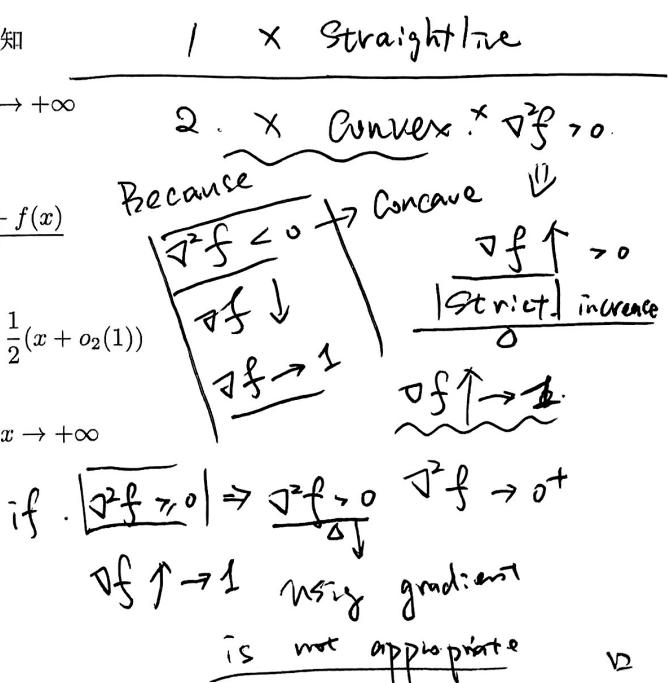
$$f(x) = x + o(1), \quad x \rightarrow +\infty$$

若 $f(x)$ 是下凸函数，那么对任意的 x ，有

$$\begin{aligned} f\left(\frac{x_0 + x}{2}\right) &\leq \frac{f(x_0) + f(x)}{2} \\ \text{Since } x \rightarrow +\infty \text{ we have} \\ \frac{x_0 + x}{2} + o_1(1) &\leq \frac{1}{2}f(x_0) + \frac{1}{2}(x + o_2(1)) \end{aligned}$$

整理得

$$\frac{x_0}{2} \leq \frac{1}{2}f(x_0) + o(1), \quad x \rightarrow +\infty$$



取极限得

$$\begin{cases} f(x) \leq x \\ f(x) = x + o(1) \quad x \rightarrow +\infty \end{cases}$$

$$\frac{x_0}{2} \leq \frac{1}{2}f(x_0), \quad x \rightarrow +\infty \Rightarrow \frac{x_0}{2} \leq \frac{f(x_0)}{2}$$

这与 $x_0 \geq 0, f(x_0) < 0$ 矛盾。所以 $f(x)$ 不可能为下凸函数。

(2) 分以下三种情况分别讨论: (a) $x_0 > 0$ (b) $x_0 < 0$ (c) $x_0 = 0$

(a) 若 $\exists x_0 > 0$ 使得 $f(x_0) = 0$, 那么 $f(x_0) = 0 < x_0$, 此式表明: 当 $x = x_0$ 时, 曲线 $y = f(x)$ 位于渐近线 $y = x$ 下方, 且竖直距离为

$$\epsilon_0 = x_0 - f(x_0) > 0$$

因为 $\lim_{x \rightarrow +\infty} (f(x) - x) = 0$, 故 $\exists \Delta > 0$, 当 $x_1 > \Delta$ 时

$$|f(x_1) - x_1| < \epsilon_0 = x_0 - f(x_0)$$

于是

$$x_1 - f(x_1) \leq |f(x_1) - x_1| < \epsilon_0 = x_0 - f(x_0)$$

即

$$x_1 - x_0 \leq f(x_1) - f(x_0)$$

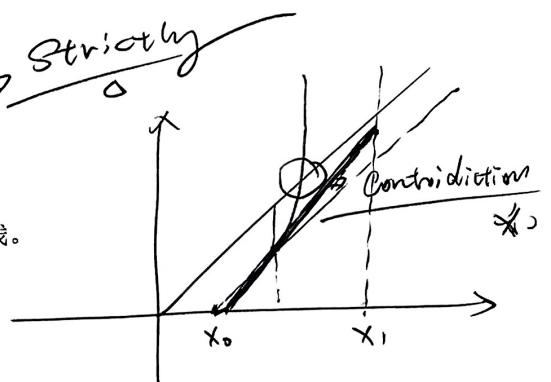
又因为 $f(x)$ 可导, 故由 Lagrange 中值定理, $\exists \xi_1 \in (x_0, x_1)$ 使得

$$(f'_1, f'_{\xi_1}) \downarrow f'(\xi_1) = \frac{f(x_1) - f(x_0)}{x_1 - x_0} \quad \text{①}$$

记 $k = f'(\xi_1)$, 则 $y = k(x - \xi_1) + f(\xi_1)$ 是 $f(x)$ 在 $x = \xi_1$ 处的切线。

根据凸函数的性质, 曲线 $f(x)$ 总在切线上方, 故

$$\underline{f(x) \geq k(x - \xi_1) + f(\xi_1)}$$



由此得

$$f(x) - x \geq (k - 1)x - k\xi_1 + f(\xi_1) \rightarrow +\infty \quad (x \rightarrow +\infty)$$

与 $y = |x|$ 是 $f(x)$ 的渐近线矛盾。

$$\underline{f(x) \leq x} \quad \Leftarrow \quad \text{if } f'(ξ₁) < 1$$

(b) 若 $\exists x_0 < 0$ 使得 $f(x_0) = 0$, 类似可证当 $x \rightarrow -\infty$ 时产生矛盾。

(c) 若 $x_0 = 0$, 即 $f(0) = 0$, 首先判断 $f(x) = |x|$ 不会发生, 否则在 $x = 0$ 处不可导。

若 $f(x) \neq |x|$, 则 $\exists x_1 \neq 0$ 使得 $f(x_1) < x_1$ 或 $f(x_1) > x_1$

(c₁) 若 $f(x_1) < x_1$ 则与 (a) 同理

(c₂) 若 $f(x_1) > x_1 > 0$, 因为 $f(0) = 0$, 可在 $[0, x_1]$ 上对 $f(x)$ 使用 Lagrange 中值定理, 找出 ξ_1 使 $f'(\xi_1) > 1$, 同上导出矛盾。

