**Lectures 15-16**

# Gradient Descent and Convexity

**DSC 40A, Fall 2025**

## Agenda

- Minimizing functions using gradient descent.
- Convexity.
- More examples.
  - Huber loss.
  - Gradient descent with multiple variables.

HW4 - no slp day
solution released on Sat.

HW3 grades by tomorrow

Midterm does <u>not</u> include:
- center and spread
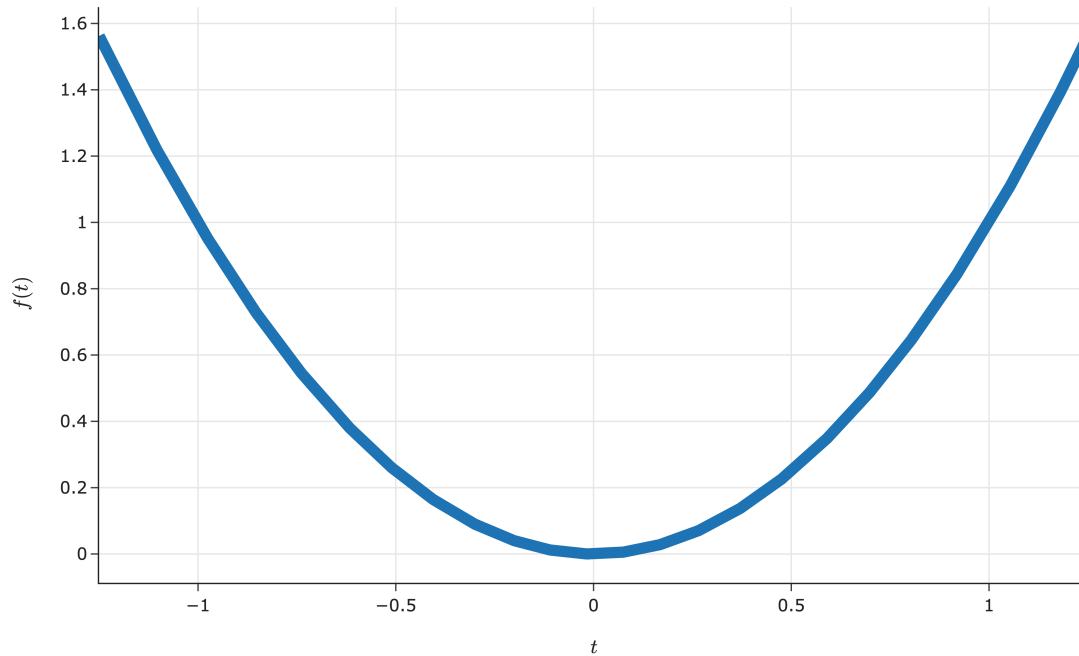  (mean absolute deviation)
- gradient descent

# Question 🤔

Answer at q.dsc40a.com

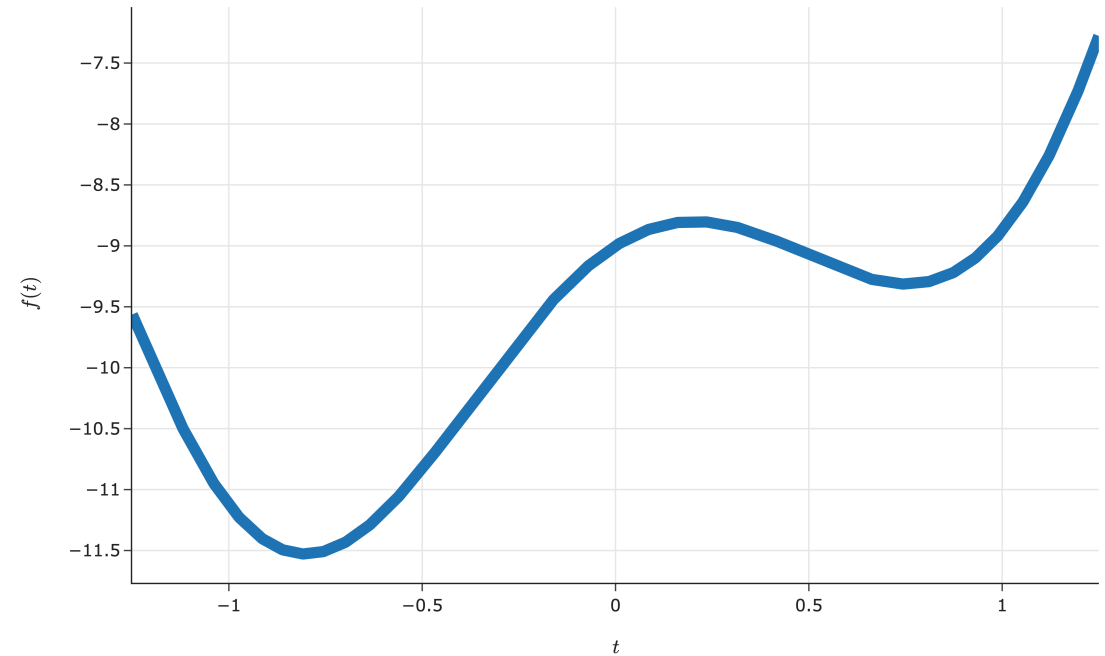**Remember, you can always ask questions at q.dsc40a.com!**

If the direct link doesn't work, click the "🤔 Lecture Questions"

link in the top right corner of dsc40a.com.

# When is gradient descent guaranteed to work?

# Convex functions



A **convex** function ✅



A **non-convex** function ❌

# Convexity

- A function $f$ is **convex** if, for **every** $a, b$ in the domain of $f$, the line segment between:

$$(a, f(a)) \text{ and } (b, f(b))$$
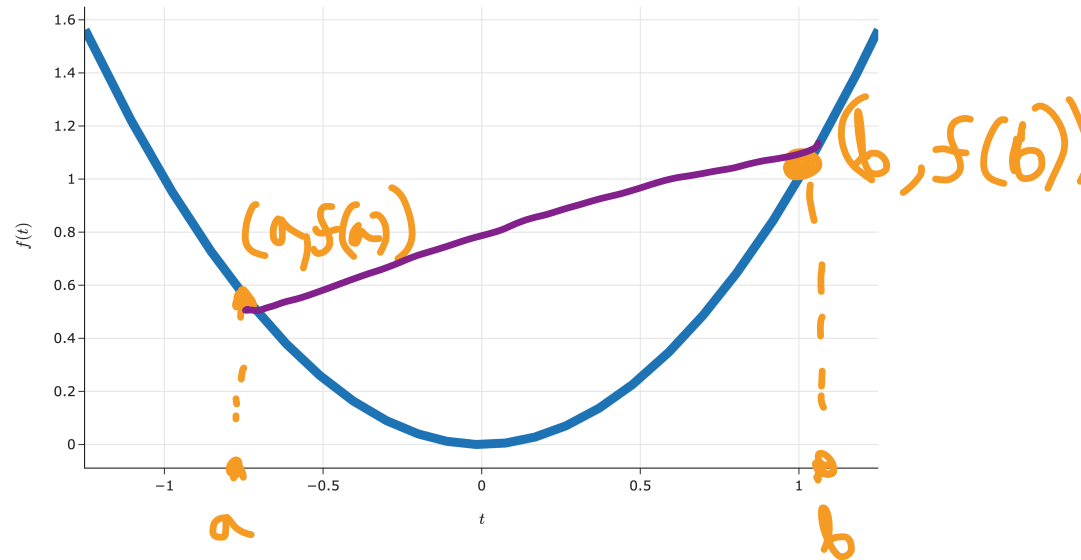
does not go below the plot of $f$.



A **convex** function ✅

# Convexity

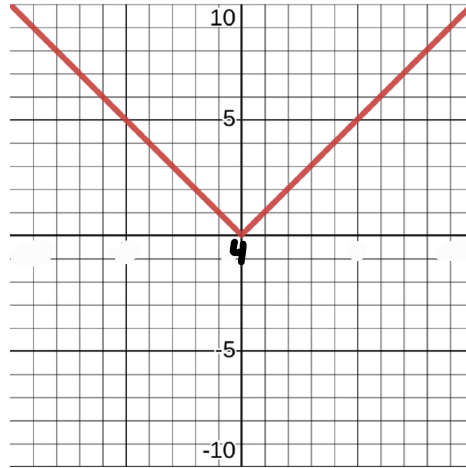- A function $f$ is **convex** if, for **every** $a, b$ in the domain of $f$, the line segment between:

$$(a, f(a)) \text{ and } (b, f(b))$$

does not go below the plot of $f$.



line is below the function
⟹ not Convex

A **non-convex** function ✗

21

# Formal definition of convexity

- A function $f : \mathbb{R} \to \mathbb{R}$ is **convex** if, for **every** $a, b$ in the domain of $f$, and for every $t \in [0, 1]$:

$$\boxed{(1-t)f(a) + tf(b) \geq f((1-t)a + tb)}$$

*(handwritten annotations:)*

plug in $t=0 \to f(a)$
plug in $t=1 \to f(b)$
plug in $t=\frac{1}{2}$, $\frac{1}{2}f(a) + \frac{1}{2}f(b)$

plug in $t=0$   $f(a)$
plug in $t=1$   $f(b)$
plug in $t=\frac{1}{2}$   $f(\frac{1}{2}a + \frac{1}{2}b)$

line between $(a, f(a)), (b, f(b))$

function between $(a, f(a)), (b, f(b))$

- A function is nonconvex if it is not convex.

- This is a formal way of restating the definition from the previous slide.

*(handwritten on right:)*

If $0 \leq t \leq 1$

$50 + 30t$

$\to (1-t)50 + 80t \qquad 0 \leq t \leq 1$

$t=0 \to 50$
$t=1 \to 80$
$t=\frac{1}{2} \to 25 + 40 = 65$

$(t=\frac{1}{2}, f(\frac{1}{2}a + \frac{1}{2}b))$

$(t=\frac{1}{2}, \frac{1}{2}f(a) + \frac{1}{2}f(b))$

$a$    $t=\frac{1}{2}$    $b$

line $\geq$ function

# Question 🤔

Answer at **q.dsc40a.com**

Is $f(x) = |x|$ convex?

- A. Yes
- B. No
- C. Maybe

$$f(x) = |x|$$

# Example: Prove $f(x) = |x|$ is convex / nonconvex

Reminder: Traingle inequality: $|\alpha + \beta| \leq |\alpha| + |\beta|$

$$(1-t)\, f(a) + t\, f(b) \geqslant f\big((1-t)a + t\,b\big) \qquad \text{for all } 0 \leq t \leq 1$$

$$\boxed{(1-t)|a| + t|b|} \geq |(1-t)a + tb|$$

$$\underbrace{|(1-t)a + tb|}_{\text{function}} \leq \underbrace{|(1-t)a| + |tb|}_{\text{triangle inequality}} \underset{0 \leq t \leq 1}{=} \boxed{\underbrace{(1-t)|a| + t|b|}_{\text{line segment}}}$$

24

# Question 🤔

Which of these functions are **not** convex?

- A. $f(x) = |x - 4|$.
- B. $f(x) = e^x$.
- C. $f(x) = \sqrt{x - 1}$.
- D. $f(x) = (x - 3)^{24}$.
- E. More than one of the above are non-convex.
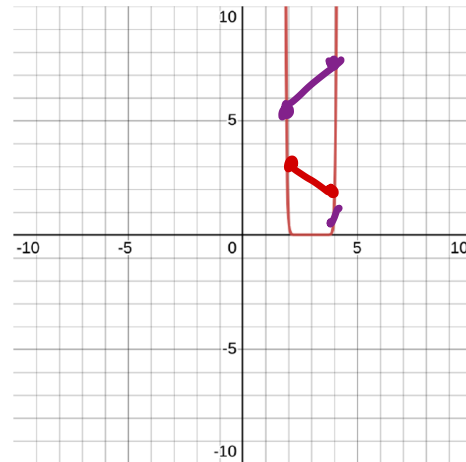
# Convex vs. concave

$|x-4|$

Convex

:)

Convex

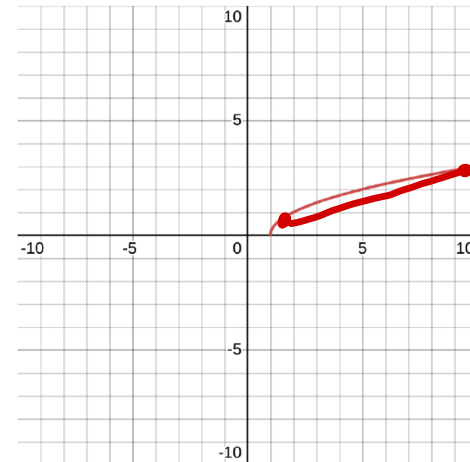:)

4

$f(x) = |x|$

$f(x) = e^x$

Convex

:)

Concave

:(

$f(x) = (x-3)^{24}$  seven

$f(x) = \sqrt{x-1}$

26

# Concave functions

- A **concave** function is the **negative** of a convex function.

# Second derivative test for convexity

- If $f(t)$ is a function of a single variable and is **twice** differentiable, then $f(t)$ is
  - convex **if and only if**:

    $\Longleftrightarrow$

    $$\frac{d^2 f}{dt^2}(t) \geq 0, \quad \forall t$$

    for all t

  - concave **if and only if**:

    $$\frac{d^2 f}{dt^2}(t) \leq 0, \quad \forall t$$

- Example: $f(x) = x^4$ is convex.

    $f'(x) = 4x^3$

    $f''(x) = 12x^2 \geq 0 \quad \forall x \implies$ convex

28

# Why does convexity matter?

- Convex functions are (relatively) easy to minimize with gradient descent.

- **Theorem**: If $f(t)$ is convex and differentiable, then gradient descent converges to a **global minimum** of $f$, as long as the step size is small enough.

- **Why?**

  ○ Gradient descent converges when the derivative is 0.

  ○ For convex functions, the derivative is 0 only at one place – the global minimum.

  ○ In other words, if $f$ is convex, gradient descent won't get "stuck" and terminate in places that aren't global minimums (local minimums, saddle points, etc.).

Step size is too large

# Nonconvex functions and gradient descent

- We say a function is **nonconvex** if it does not meet the criteria for convexity.

- Nonconvex functions are (relatively) difficult to minimize.

- Gradient descent **might** still work, but it's not guaranteed to find a global minimum.

  - We saw this at the start of the lecture, when trying to minimize $f(t) = 5t^4 - t^3 - 5t^2 + 2t - 9$.

global min.

← local min

30

## Choosing a step size in practice

- In practice, choosing a step size involves a lot of trial-and-error.

- In this class, we've only touched on "constant" step sizes, i.e. where $\alpha$ is a constant.

$$t_{i+1} = t_i - \alpha \frac{df}{dt}(t_i)$$

- **Remember**: $\alpha$ is the "step size", but the amount that our guess for $t$ changes is $\alpha \frac{df}{dt}(t_i)$, not just $\alpha$.

- In future courses, you'll learn about "decaying" step sizes, where the value of $\alpha$ decreases as the number of iterations increases.

  - Intuition: take much bigger steps at the start, and smaller steps as you progress, as you're likely getting closer to the minimum.

# More examples

$$H(x) = h$$

# Example: Huber loss and the constant model

- First, we learned about squared loss,
  $$L_{\text{sq}}(y_i, H(x_i)) = (y_i - H(x_i))^2.$$
  pro: differentiable, easy to minimize
  con: sensitive to outliers

- Then, we learned about absolute loss,
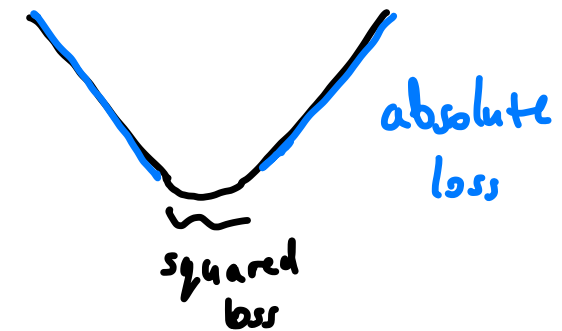  $$L_{\text{abs}}(y_i, H(x_i)) = |y_i - H(x_i)|.$$
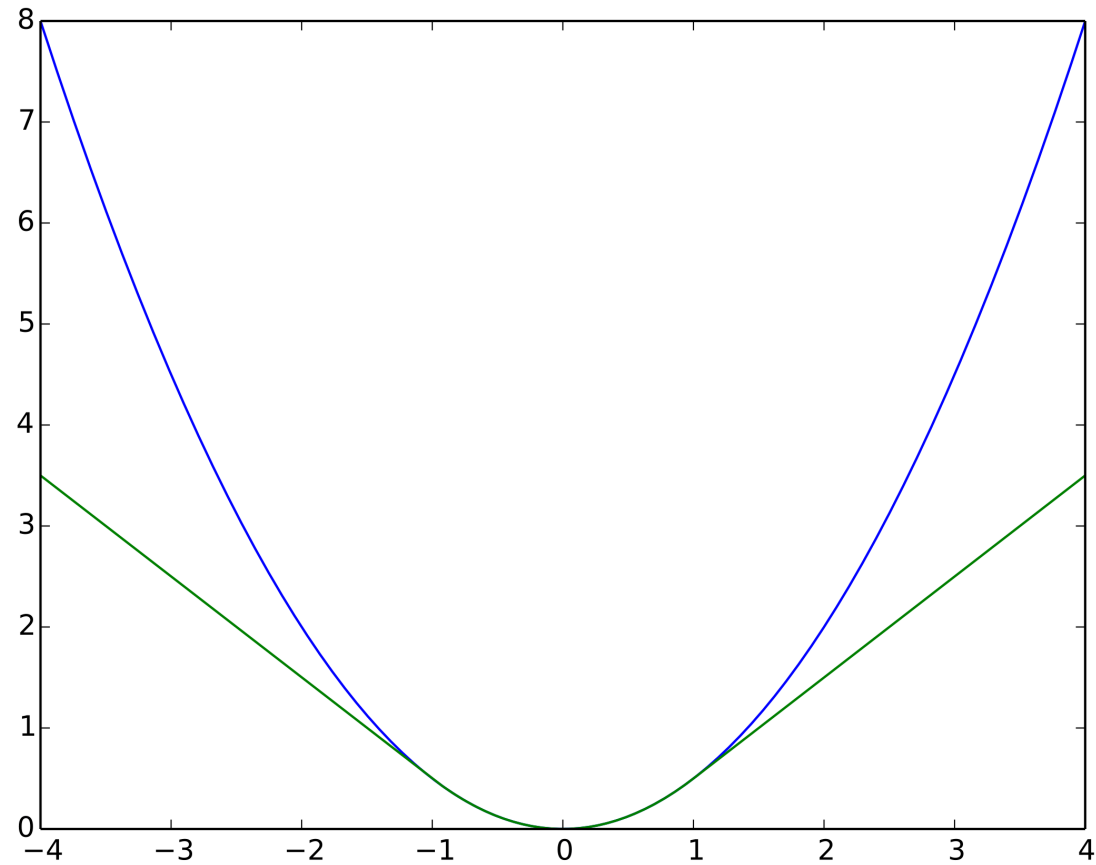  pro: robust to outliers

  con: not differentiable

  absolute loss

  squared loss

- Let's look at a new loss function, **Huber loss**:

$$L_{\text{huber}}(y_i, H(x_i)) = \begin{cases} \frac{1}{2}(y_i - H(x_i))^2 & \text{if } |y_i - H(x_i)| \leq \delta \longrightarrow \text{hyperparam} \\ \delta \cdot (|y_i - H(x_i)| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

**Squared** loss in blue, **Huber** loss in green.
Note that both loss functions are convex!

# Minimizing average Huber loss for the constant model

- For the constant model, $H(x) = h$:

$$L_{\text{huber}}(y_i, h) = \begin{cases} \frac{1}{2}(y_i - h)^2 & \text{if } |y_i - h| \leq \delta \\ \delta \cdot (|y_i - h| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

$$\implies \frac{\partial L}{\partial h}(h) = \begin{cases} -(y_i - h) & \text{if } |y_i - h| \leq \delta \\ -\delta \cdot \text{sign}(y_i - h) & \text{otherwise} \end{cases}$$

$$\text{sign}(x) = \begin{cases} +1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases}$$

- So, the **derivative** of empirical risk is:

$$\frac{dR_{\text{huber}}}{dh}(h) = \frac{1}{n} \sum_{i=1}^{n} \begin{cases} -(y_i - h) & \text{if } |y_i - h| \leq \delta \\ -\delta \cdot \text{sign}(y_i - h) & \text{otherwise} \end{cases}$$

- It's **impossible** to set $\frac{dR_{\text{huber}}}{dh}(h) = 0$ and solve by hand: we need gradient descent!

$$R_{sq} = \frac{1}{n} \sum_i \ell_{\text{huber}}(y_i, h)$$

35

Let's try this out in practice! Follow along in this notebook.