

EECS 245, Spring 2026

LEC 10

Gradients and
Gradient Descent

→ Read: Ch. 8

Agenda

Ch 8.1-8.4
all in scope

- Recap: Big 3 gradient rules

- Minimizing

$$R_{sq}(\vec{w}) = \frac{1}{n} \|\vec{y} - X\vec{w}\|^2$$

using gradients

- Using gradients to minimize functions algorithmically - gradient descent

- If time: convexity -
not in scope for MT2,
in scope for Final

Announcements

- Deadlines:

- HW 7 today

- HW 8 Sunday, no slip days

- Lab 9 Monday

- Midterm 2 on Tuesday
from 1-3PM! 2 ^{new} conceptual Q's

- No in-person lecture
next Thursday - videos
will be posted

Practice Conceptual Questions

It's often said that the best way to learn something is by teaching it. In IA interviews for this class, we usually ask conceptual questions like the ones below. It's a good idea to study for Midterm 2 by trying to explain the questions below to your peers and having them ask questions of your explanation.

- 1 Explain where the normal equations came from.
- 2 How does finding the line of best fit have anything to do with the normal equations?
- 3 Prove that the null space is a subspace.
- 4 Prove that every element in the column space is orthogonal to every element in the null space of A^T .
- 5 What does the rank-nullity theorem say?

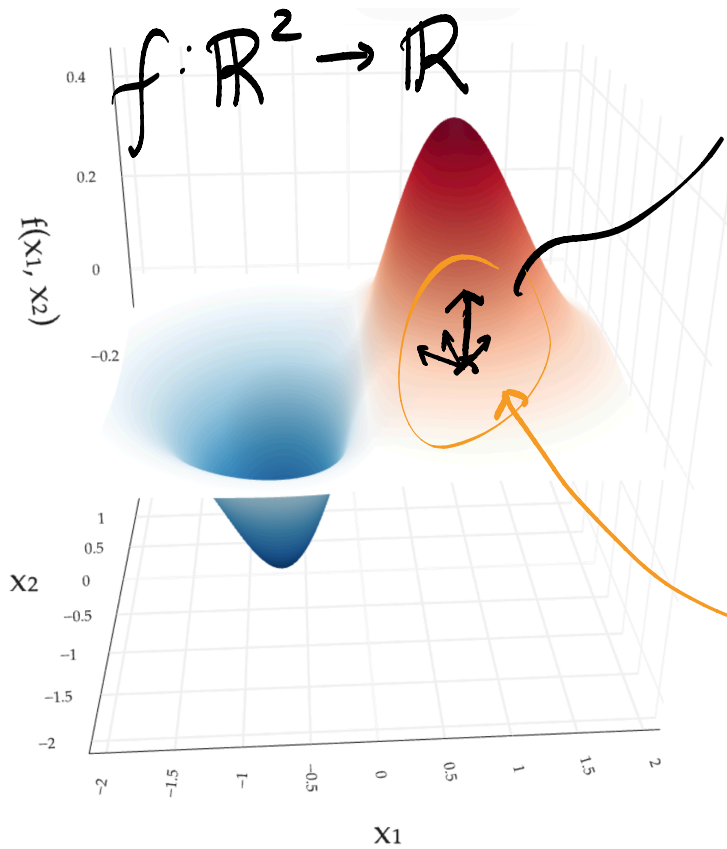
Recap: Gradients

$$f: \mathbb{R}^d \rightarrow \mathbb{R}$$

"vector-to-scalar"

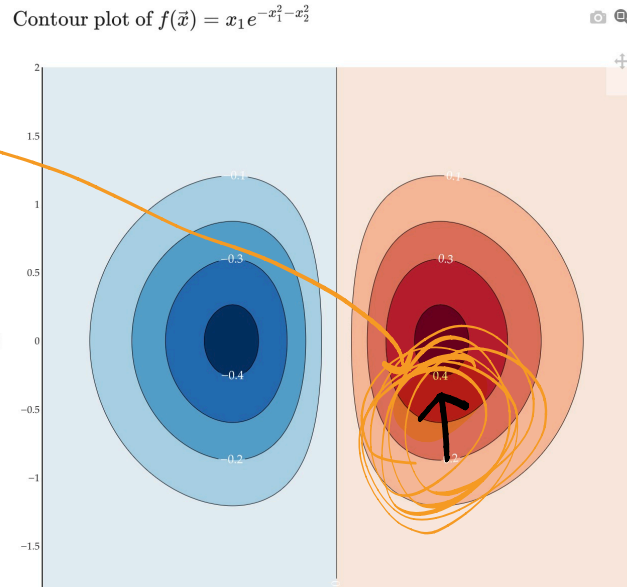
e.g. $f(\vec{x}) = (x_1 + x_2)^2 - 3\cos x_1$ $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$\nabla f(\vec{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2(x_1 + x_2) + 3\sin x_1 \\ 2(x_1 + x_2) \end{bmatrix}$$



function increases in lots of directions, but increases most quickly in the direction of the gradient,

Contour plot of $f(\vec{x}) = x_1 e^{-x_1^2 - x_2^2}$



Ultimate goal: minimize mean squared error

$$R_{sq}(\vec{w}) = \frac{1}{n} \|\underbrace{\vec{y} - X\vec{w}}_{\text{error vector}}\|^2$$

Why? The \vec{w}^* that minimizes $R_{sq}(\vec{w})$ will tell us

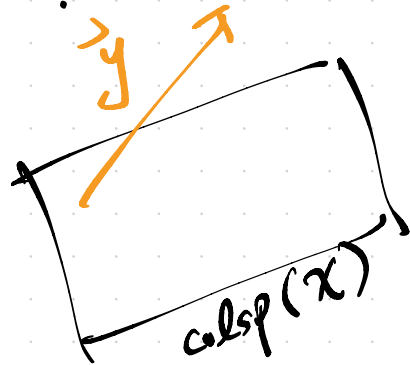
how to make the best predictions!

(\vec{w}^* - optimal model parameters)

We know what the answer should be!

If $X^T X$ invertible,

$$\vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$



Recap: Big 3 gradient rules

$$\textcircled{1} f(\vec{x}) = \vec{a}^T \vec{x} = \vec{a} \cdot \vec{x} \quad \vec{x}, \vec{a} \in \mathbb{R}^n$$

$$\nabla f(\vec{x}) = \vec{a}$$

$$\textcircled{2} f(\vec{x}) = \|\vec{x}\|^2 \quad \vec{x} \in \mathbb{R}^n$$

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

"quadratic form"

$$\textcircled{3} f(\vec{x}) = \vec{x}^T A \vec{x} \quad \vec{x} \in \mathbb{R}^n, A: n \times n \text{ matrix}$$

$$\nabla f(\vec{x}) = (A + A^T) \vec{x}$$

e.g.

$$f(\vec{x}) = \|\vec{x}\|$$

$$\vec{x} \xrightarrow{\text{vector to scalar } g} \|\vec{x}\|^2 \xrightarrow{\text{scalar to scalar } h} \sqrt{\|\vec{x}\|^2} = \|\vec{x}\|$$

We would like to use the fact that

$$\nabla(\|\vec{x}\|^2) = 2\vec{x}$$

$$f(\vec{x}) = h(g(\vec{x}))$$

Recall, $\frac{d}{dx} \sqrt{x} = \frac{1}{2\sqrt{x}}$

$$\nabla f(\vec{x}) = \frac{dh}{dx}(g(\vec{x})) \nabla g(\vec{x}) = h'(g(\vec{x})) \nabla g(\vec{x})$$

$$\text{Here, } \nabla \|\vec{x}\| = \nabla(\sqrt{\|\vec{x}\|^2}) = \frac{1}{2\sqrt{g(\vec{x})}} \nabla g(\vec{x}) = \frac{\vec{x}}{\|\vec{x}\|}$$

$= \frac{1}{\|\vec{x}\|} \cdot 2\vec{x}$

$$f(\vec{x}) = \|\vec{x}\|^p$$

$p \in \mathbb{R}$

Chapter 8.2

$$f(\vec{x}) = h(g(\vec{x}))$$

$$h(x) = x^p \quad g(\vec{x}) = \|\vec{x}\|$$

$$\nabla f(\vec{x}) = \frac{dh}{dx}(g(\vec{x})) \nabla g(\vec{x})$$

$$= p g(\vec{x})^{p-1} \frac{\vec{x}}{\|\vec{x}\|} = \frac{p \|\vec{x}\|^{p-1} \vec{x}}{\|\vec{x}\|} = p \|\vec{x}\|^{p-2} \vec{x}$$

$$R_{sq}(\vec{w}) = \frac{1}{n} \|\vec{y} - X\vec{w}\|^2$$

want to find

$$\nabla R_{sq}(\vec{w})$$

$$= \frac{1}{n} (\vec{y} - X\vec{w}) \cdot (\vec{y} - X\vec{w})$$

$$= \frac{1}{n} (\vec{y} \cdot \vec{y} - \vec{y} \cdot (X\vec{w}) - (X\vec{w}) \cdot \vec{y} + (X\vec{w}) \cdot (X\vec{w}))$$

$$= \frac{1}{n} (\vec{y} \cdot \vec{y} - 2(X\vec{w}) \cdot \vec{y} + (X\vec{w})^T (X\vec{w}))$$

$$= \frac{1}{n} (\vec{y} \cdot \vec{y} - 2\vec{w}^T (X^T \vec{y}) + \vec{w}^T X^T X \vec{w})$$

X : $n \times d$ matrix

$$\vec{y} \in \mathbb{R}^n$$

$$\vec{w} \in \mathbb{R}^d$$

$$\vec{u} \cdot \vec{v} = \vec{u}^T \vec{v}$$

Aside:

$$\begin{aligned} & \underbrace{\vec{u}}_{(\chi \vec{w})} \cdot \underbrace{\vec{v}}_{\vec{y}} \\ &= \vec{w}^T \underbrace{\chi^T}_{\vec{y}} \end{aligned}$$

$$\text{so, } (\chi \vec{w}) \cdot \vec{y} = \vec{w} \cdot (\chi^T \vec{y})$$

$$R_{sq}(\vec{w}) = \frac{1}{n} \left(\vec{y} \cdot \vec{y} - \underbrace{2\vec{w}^T (\overset{\vec{a}}{X^T \vec{y}})} + \vec{w}^T \boxed{X^T X} \vec{w} \right)$$

$$\nabla R_{sq}(\vec{w}) = \frac{1}{n} \left(-2X^T \vec{y} + \underbrace{\left((X^T X) + (X^T X)^T \right)}_{2X^T X} \vec{w} \right)$$

$$= -\frac{2}{n} \left(X^T \vec{y} - X^T X \vec{w} \right)$$

set
 $\nabla R_{SS}(\vec{w})$

$$-\frac{2}{n} (X^T \vec{y} - X^T X \vec{w}) = \vec{0}$$

$$X^T \vec{y} - X^T X \vec{w} = \vec{0}$$

the normal
equation!

→

$$X^T X \vec{w} = X^T \vec{y}$$

same condition

for \vec{w} we
saw before

with orthogonality and
projections!

$$X^T X \bar{w} = X^T \bar{y}$$

all \bar{w} 's that minimize

$$R_{sq}(\bar{w}) = \frac{1}{n} \|\bar{y} - X\bar{w}\|^2$$

satisfy

→ if $X^T X$ invertible, there is a unique best \bar{w}^* :

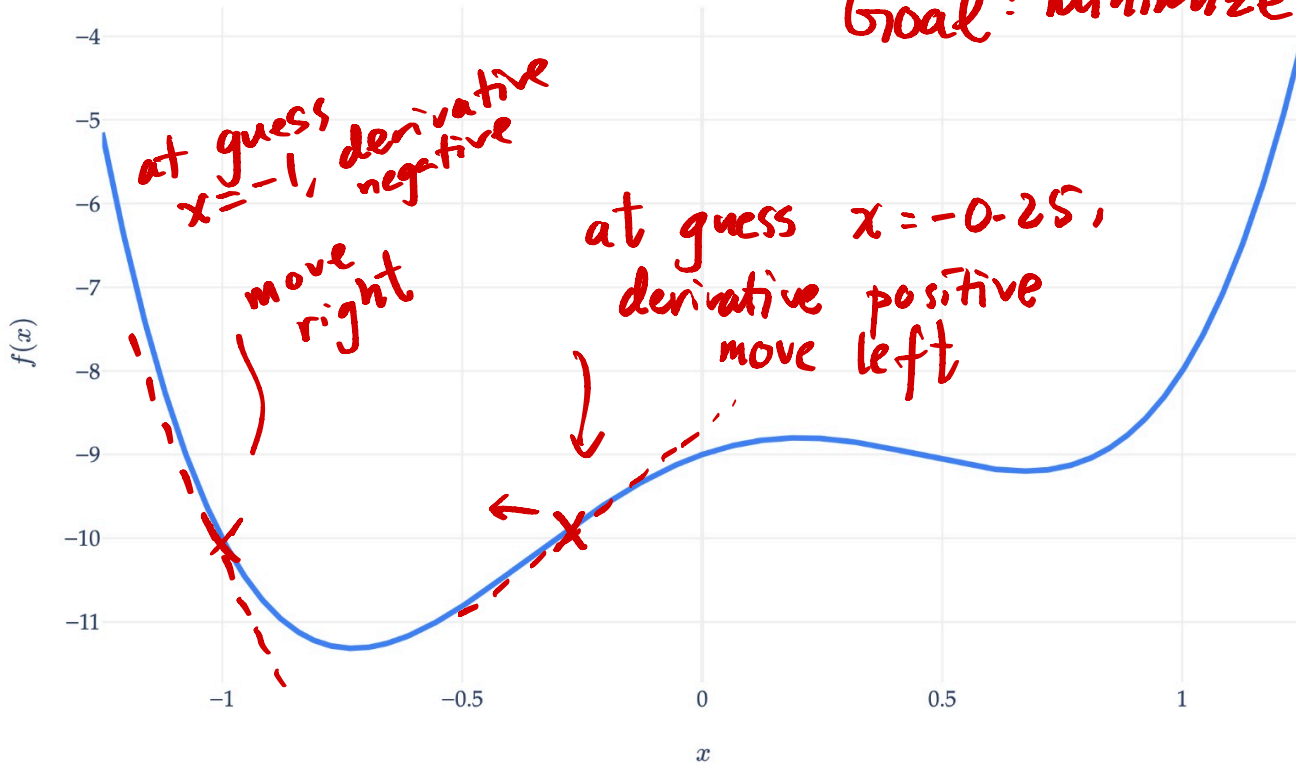
$$\bar{w}^* = (X^T X)^{-1} X^T \bar{y}$$

→ otherwise, infinitely many — see Chapter 6.4

$$f(x) = 5x^4 - x^3 - 5x^2 + 2x - 9$$

$$f'(x) = 20x^3 - 3x^2 - 10x + 2$$

Goal: minimize f



Gradient descent

Goal: Minimize $f: \mathbb{R}^d \rightarrow \mathbb{R}$, a differentiable function

First, choose initial guess, $\vec{x}^{(0)}$, and a learning rate / step size, $\alpha > 0$

update guess using update rule:

$$\vec{x}^{(t+1)} = \vec{x}^{(t)} - \alpha \underbrace{\nabla f(\vec{x}^{(t)})}_{\text{vector}}$$

Terminate when $\|\nabla f(\vec{x}^{(t)})\| \leq 0.001$ (tolerance)

Many examples

in Ch 8.3 - go play with them!

Chapter 8.3

Consider the following function.

$$f(\vec{x}) = (x_1 - 2)^2 + 2x_1 - (x_2 - 3)^2$$

1. Is $f(\vec{x})$ a quadratic form?

2. Given an initial guess of $\vec{x}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and a step size of $\alpha = \frac{1}{3}$, perform ~~two~~ ^{one} iterations of gradient descent. What is $\vec{x}^{(2)}$?

$$f(\vec{x}) = (x_1 - 2)^2 + 2x_1 - (x_2 - 3)^2$$

$$\nabla f(\vec{x}) = \begin{bmatrix} 2(x_1 - 2) + 2 \\ -2(x_2 - 3) \end{bmatrix}$$

$$\vec{x}^{(0)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \alpha = \frac{1}{3}$$

$$\begin{aligned} \vec{x}^{(1)} &= \vec{x}^{(0)} - \alpha \nabla f(\vec{x}^{(0)}) \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} 2(-2) + 2 \\ -2(-3) \end{bmatrix} = \begin{bmatrix} 2/3 \\ -2 \end{bmatrix} \end{aligned}$$

Issues with gradient descent:

numerical method

- Can get trapped at a local minimum
- step size too large \rightarrow may not converge

Gradient descent is usually used for finding optimal model parameters, i.e. for empirical risk minimization

e.g. $R_{sq}(\hat{w}) = \frac{1}{n} \|\hat{y} - X\hat{w}\|^2$

since we have a "closed-form" solution for \hat{w}^* , we don't strictly need GD here

Let's try it!

$$R_{sq}(\vec{w}) = \frac{1}{n} \|\vec{y} - X\vec{w}\|^2$$

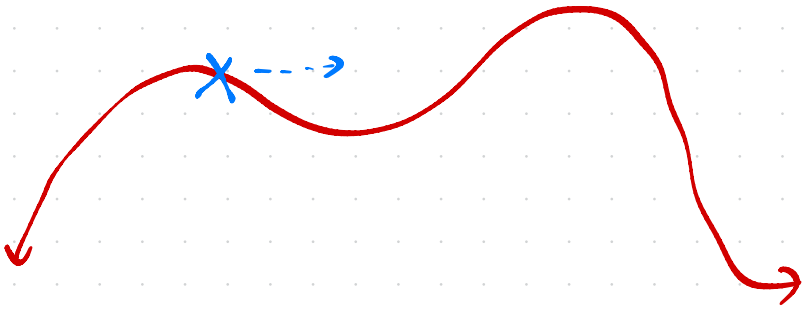
$$\nabla R_{sq}(\vec{w}) = \underbrace{-\frac{2}{n} (X^T \vec{y} - X^T X \vec{w})}_{\text{from earlier today}}$$

Start with $\vec{w}^{(0)}$

from earlier today

$$\vec{w}^{(t+1)} = \vec{w}^{(t)} - \alpha \left(-\frac{2}{n} (X^T \vec{y} - X^T X \vec{w}^{(t)}) \right)$$

Update until $\|\nabla R_{sq}(\vec{w}^{(t)})\| \leq \text{tolerance}$



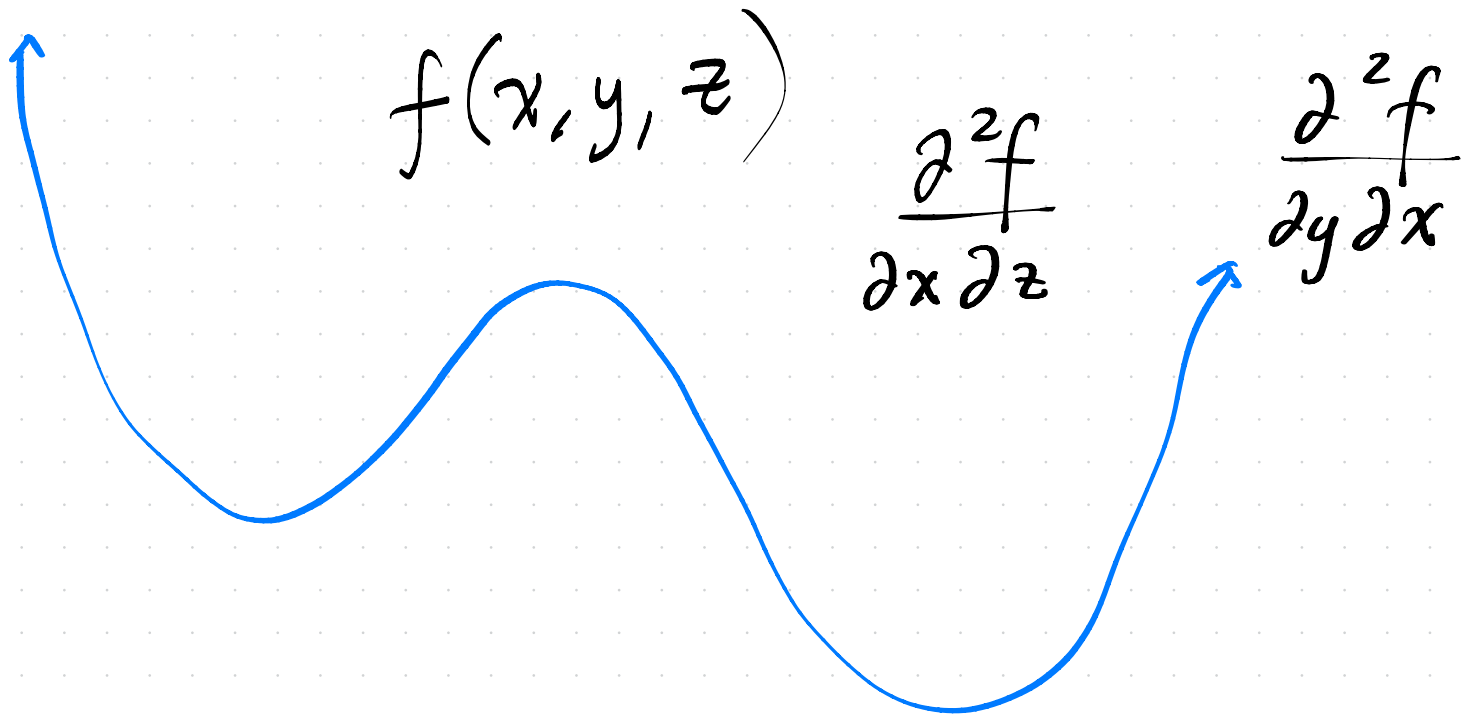
Recap: we use gradient descent to

minimize

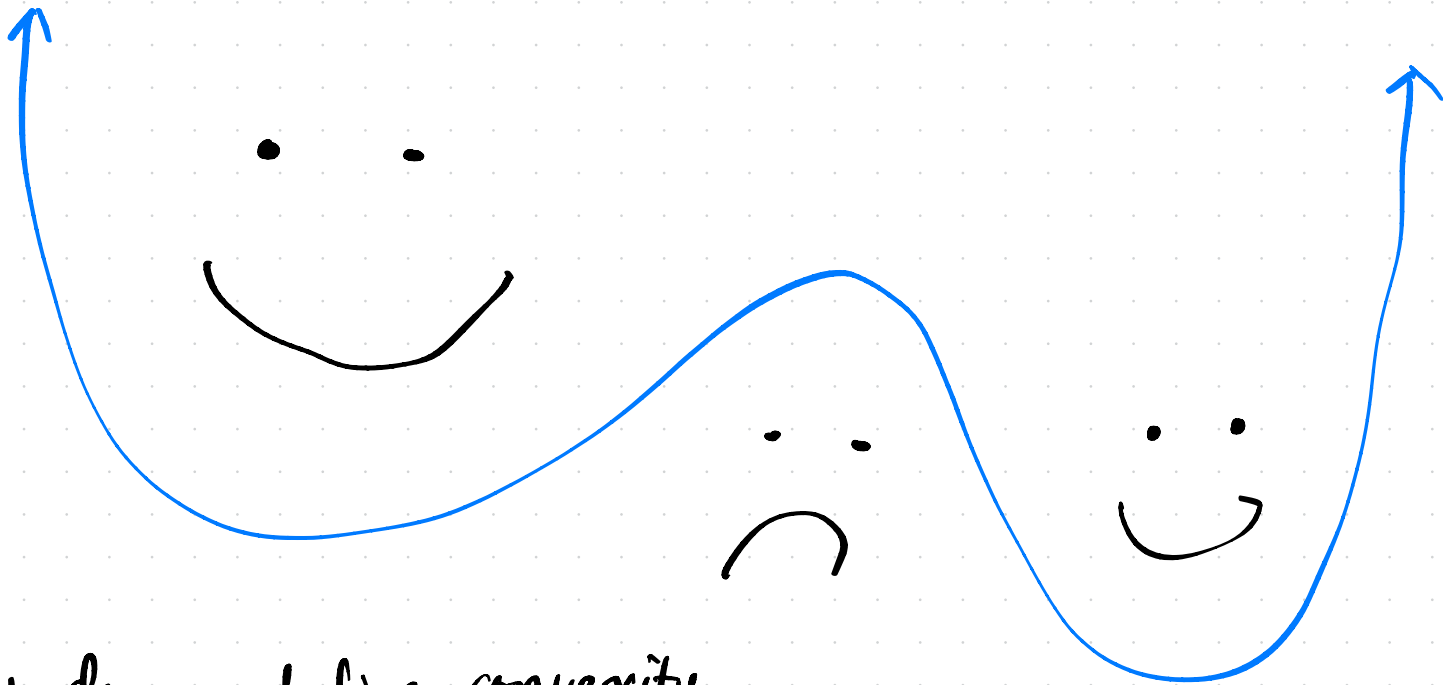
differentiable functions (vector-to-scalar
or scalar-to-scalar)

when we don't have a closed-form minimizer.

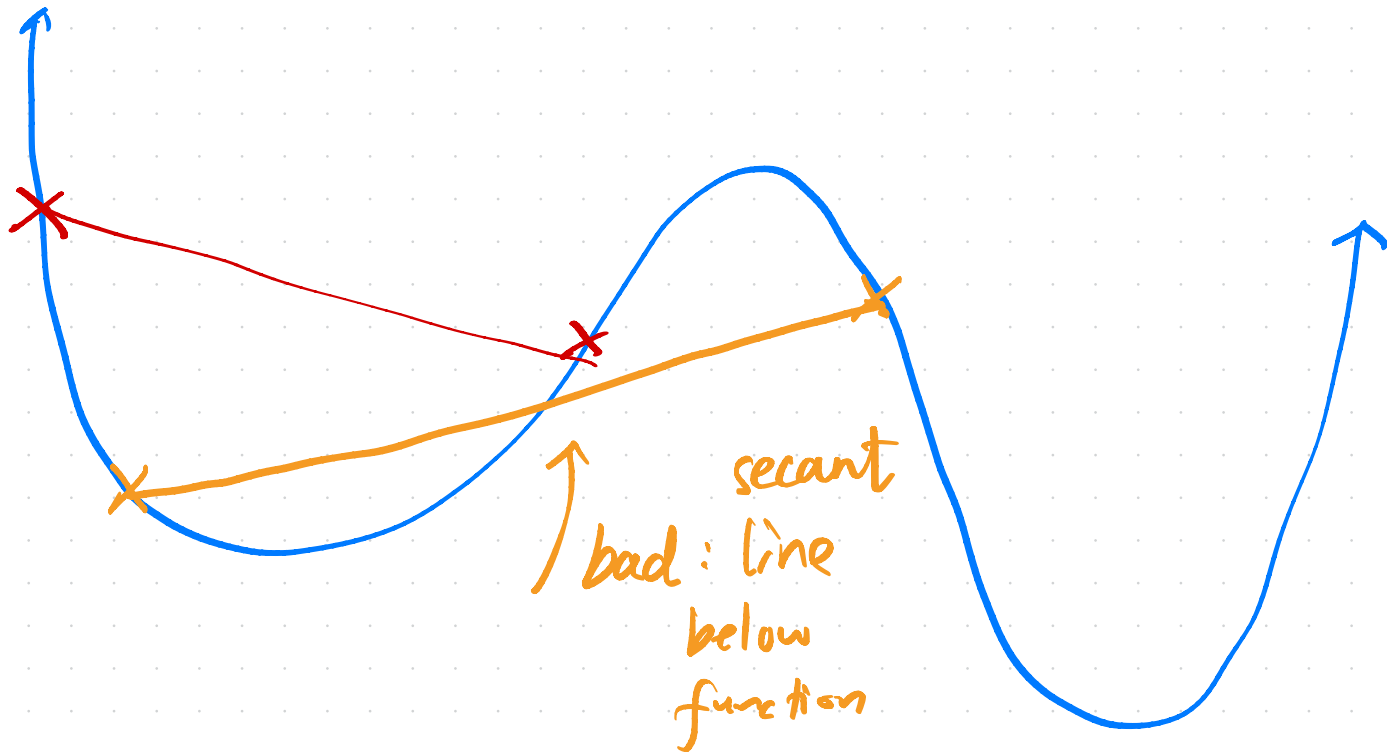
Usually, the function is average loss / empirical risk,
so the minimizer is a vector of optimal
model parameters.



→ not on MT 2,
but on final exam



How do we define convexity
without using derivatives?

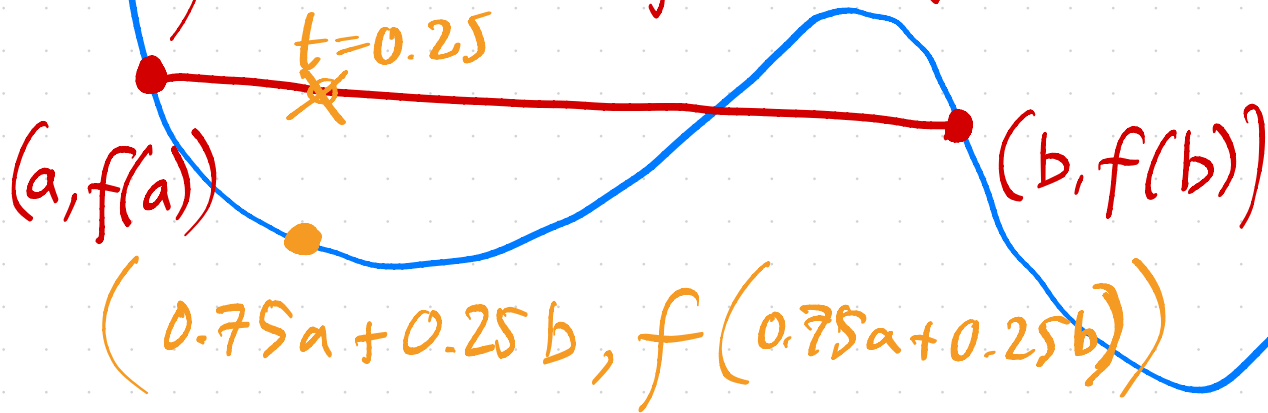


Goal: Translate

"all possible secant lines are on or above function"

to a precise mathematical expression

$$L: f(a) + t(f(b) - f(a)) \quad t \in [0, 1]$$
$$= (1-t)f(a) + tf(b)$$



Suppose $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a function.

f is convex if

$$f\left(\underbrace{(1-t)\vec{x} + t\vec{y}}_{\text{function}}\right) \leq \underbrace{(1-t)f(\vec{x}) + tf(\vec{y})}_{\text{line}}$$

for all $t \in [0, 1]$

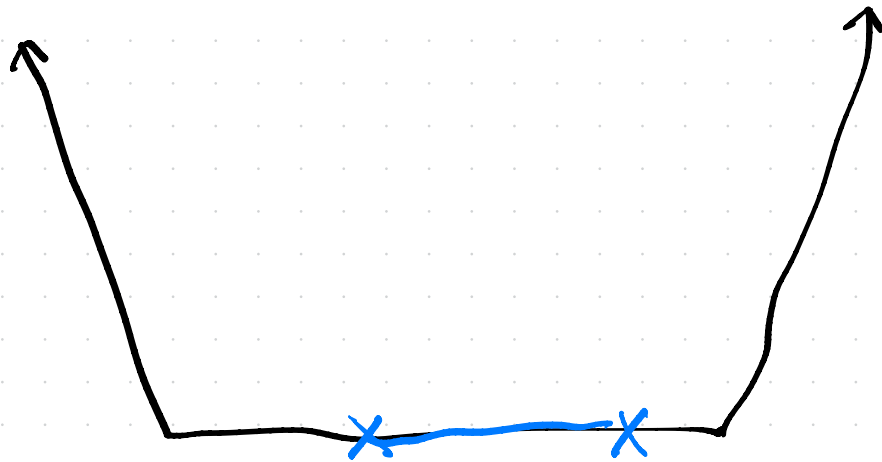
and all \vec{x}, \vec{y} in domain of f .

Why is this useful?

- If we can prove function convex,
it doesn't have any "trap"
local minimums

→ GD won't get stuck

→ see video in 8.5



convex?
yes!

strictly convex?
no!

e.g. mean squared error,

$$R_{sq}(\vec{w}) = \frac{1}{n} \|\vec{y} - X\vec{w}\|^2$$

always convex

only strictly convex if X's columns are linearly independent.