

Lectures 15-16

# Gradient Descent and Convexity

DSC 40A, Fall 2024

# Agenda

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

- HW2 Q5a
- Cheat sheet
- FAQs updated
- Regarding requests  
only through gradescope  
(Not email, not OH)

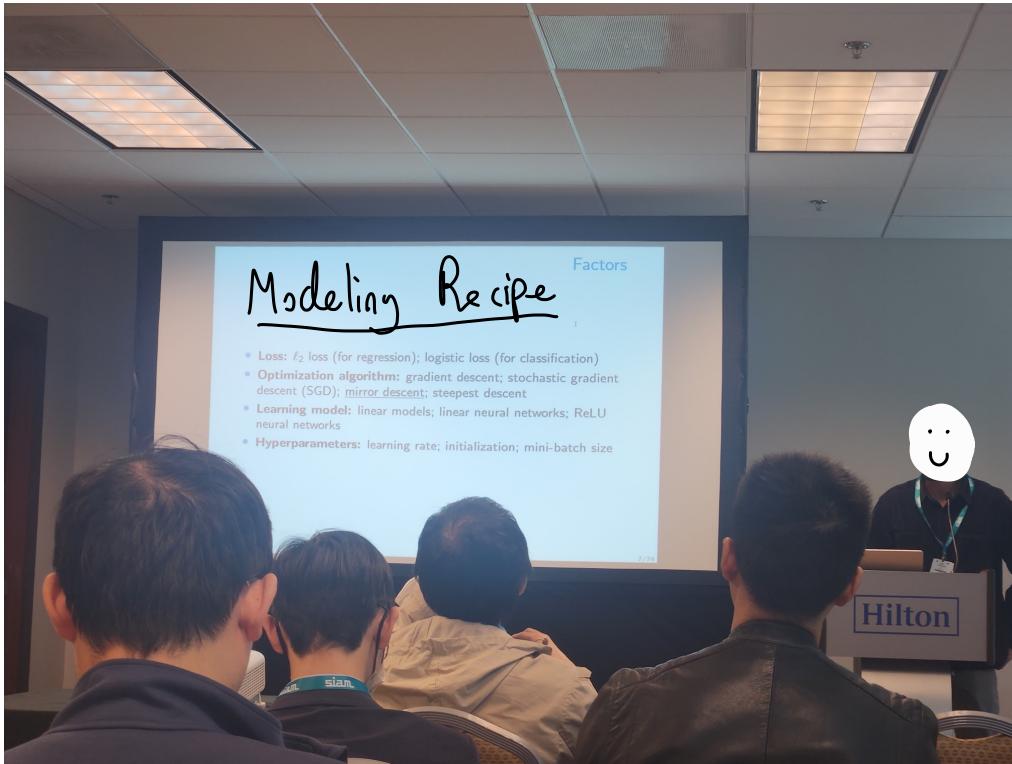
## Question 🤔

Answer at [q.dsc40a.com](https://q.dsc40a.com)

**Remember, you can always ask questions at [q.dsc40a.com!](https://q.dsc40a.com)**

If the direct link doesn't work, click the " Lecture Questions" link in the top right corner of [dsc40a.com](https://dsc40a.com).

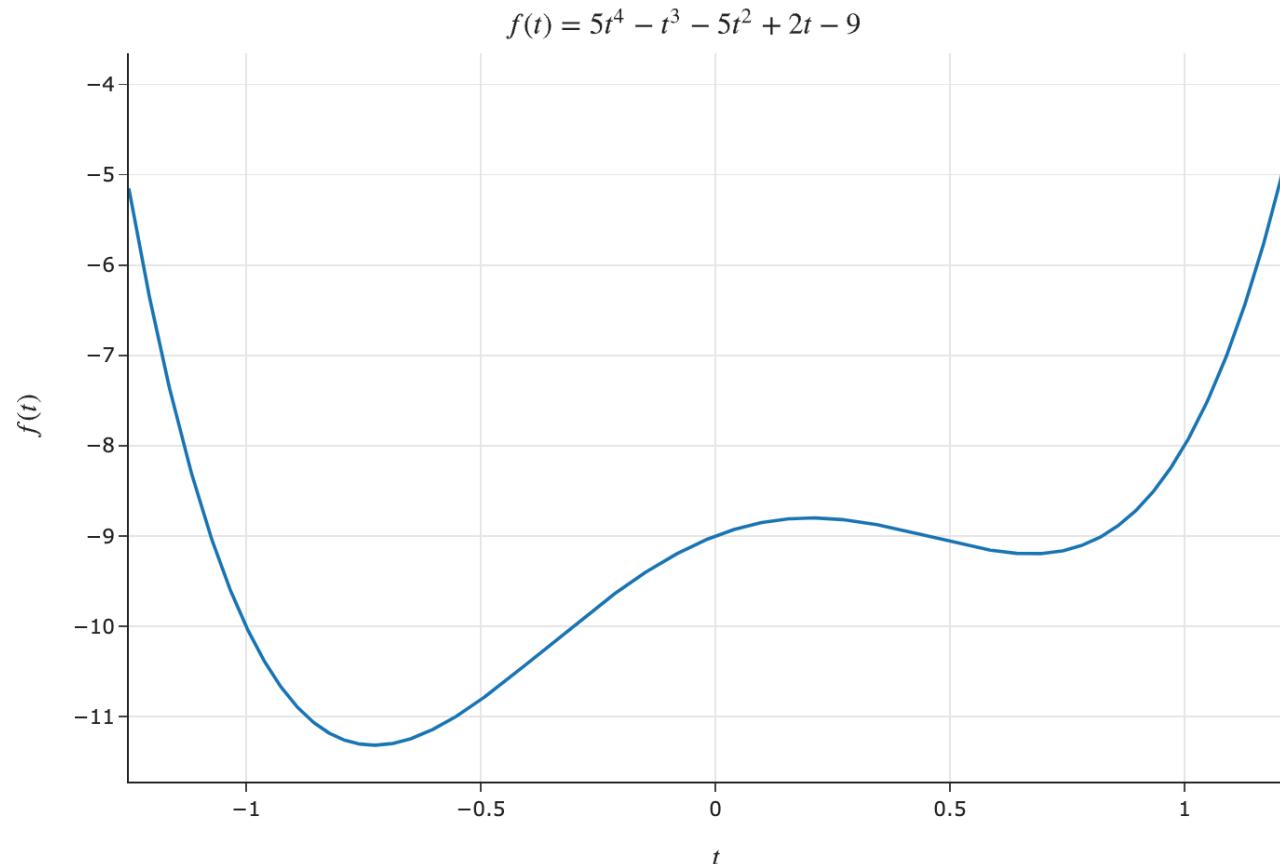
# Anecdote



# Minimizing functions using gradient descent

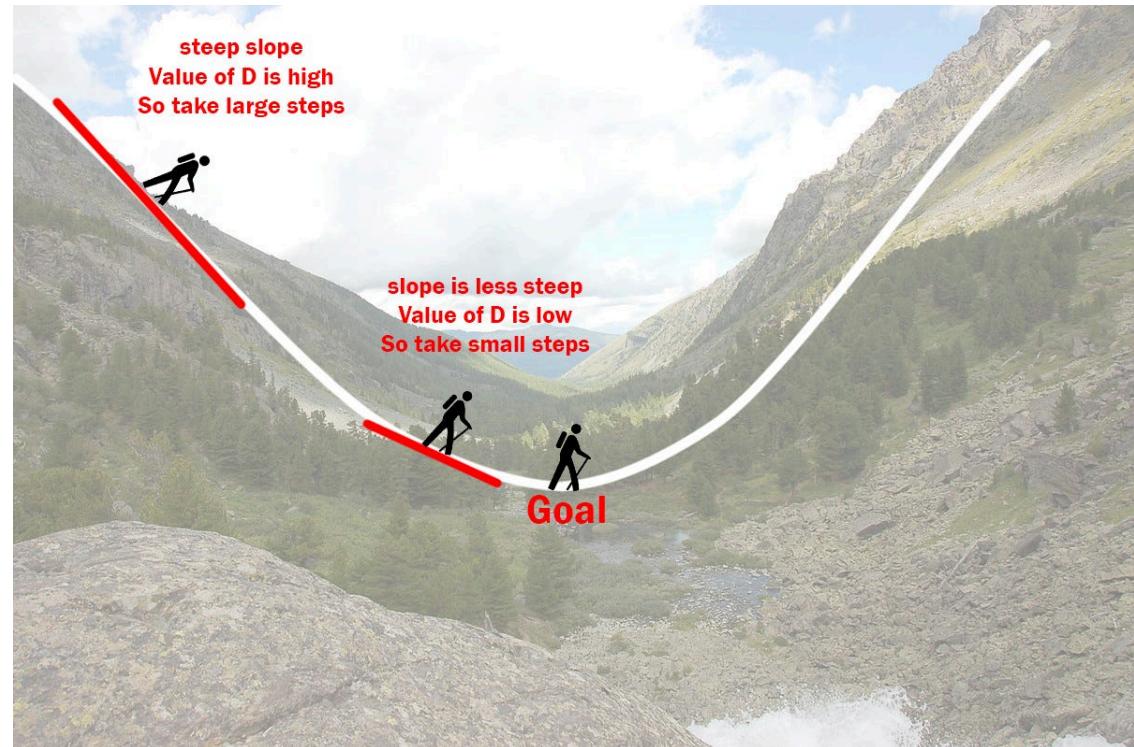
# What does the derivative of a function tell us?

- Goal: Given a differentiable function  $f(t)$ , find the input  $t^*$  that minimizes  $f(t)$ .
- What does  $\frac{d}{dt} f(t)$  mean?

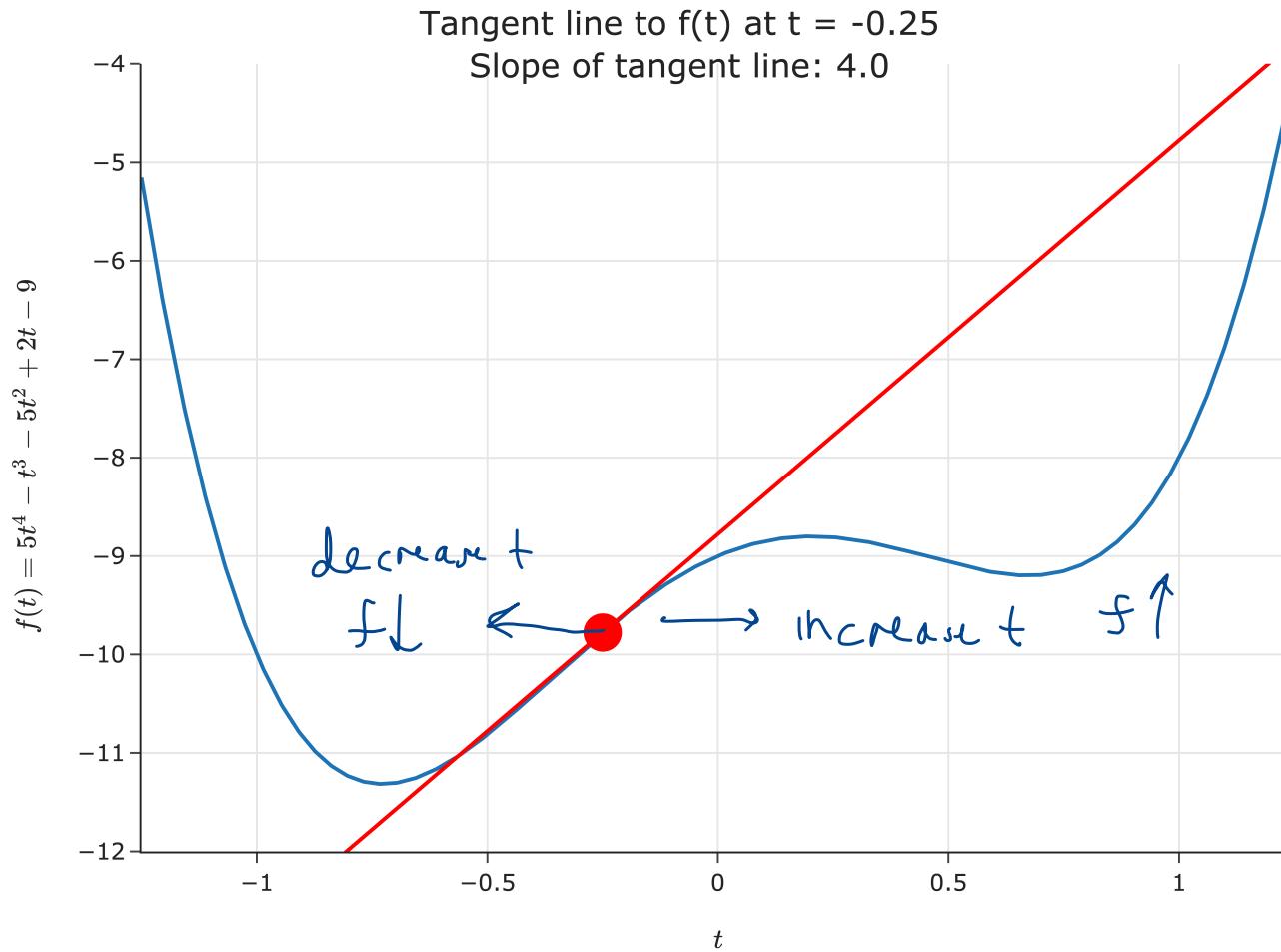


# Let's go hiking!

- Suppose you're at the top of a mountain  and need to get **to the bottom**.
- Further, suppose it's really cloudy , meaning you can only see a few feet around you.
- **How** would you get to the bottom?



# Searching for the minimum

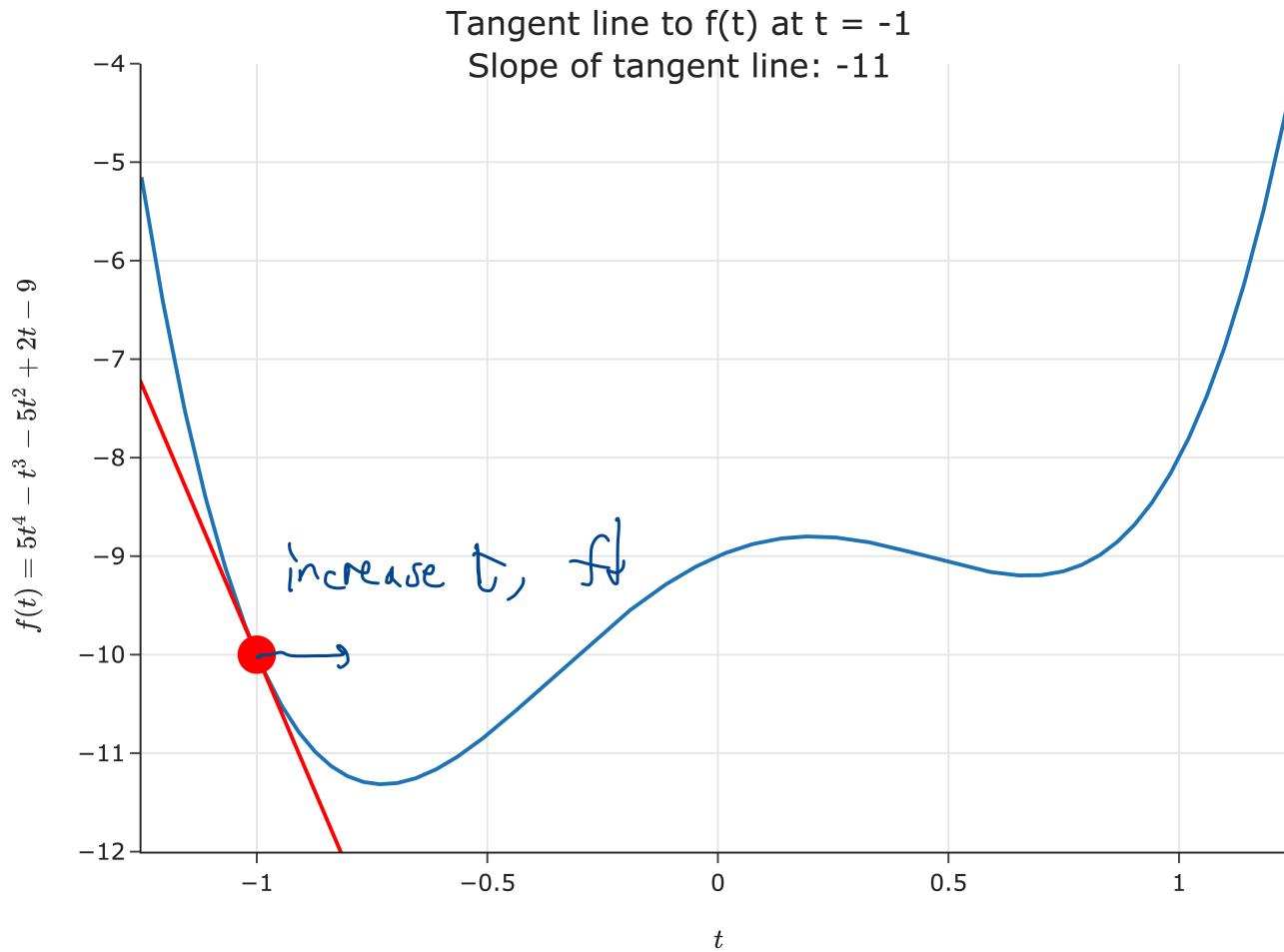


Suppose we're given an initial guess for a value of  $t$  that minimizes  $f(t)$ .

If the **slope of the tangent line at  $f(t)$  is positive ↗**:

- Increasing  $t$  increases  $f$ .
- This means the minimum must be to the **left** of the point  $(t, f(t))$ .
- Solution: Decrease  $t$  ↓.

# Searching for the minimum



Suppose we're given an initial guess for a value of  $t$  that minimizes  $f(t)$ .

If the **slope of the tangent line at  $f(t)$  is negative** :

- Increasing  $t$  decreases  $f$ .
- This means the minimum must be to the **right** of the point  $(t, f(t))$ .
- Solution: **Increase  $t$**  .

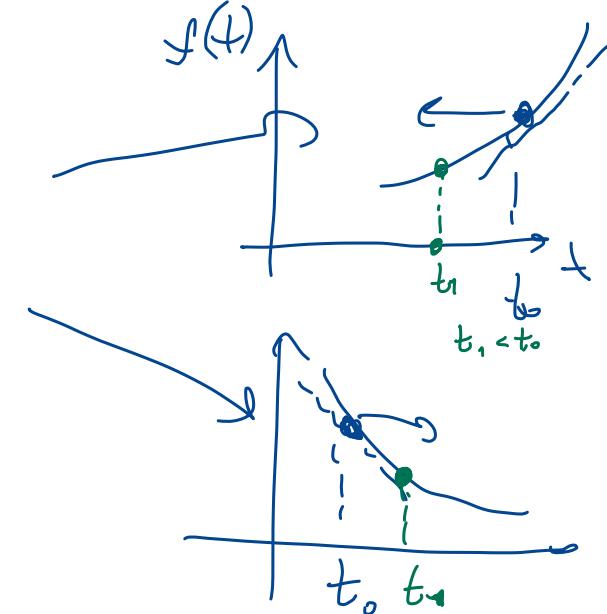
# Intuition

- To minimize  $f(t)$ , start with an initial guess  $t_0$ .
- Where do we go next?
  - If  $\frac{df}{dt}(t_0) > 0$ , decrease  $t_0$ .
  - If  $\frac{df}{dt}(t_0) < 0$ , increase  $t_0$ .
- One way to accomplish this:

$$t_1 = t_0 - \boxed{\text{?}}$$
$$t_1 = t_0 + \boxed{\text{?}}$$

$$t_1 = t_0 - \underbrace{\frac{df}{dt}(t_0)}$$

opposite direction of  
the derivative

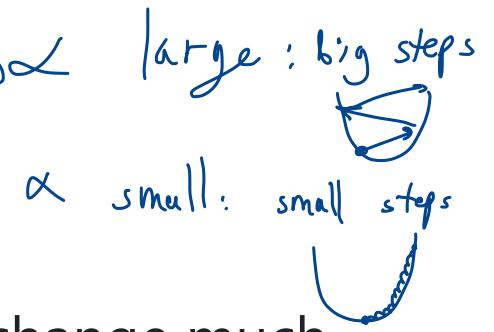


# Gradient descent

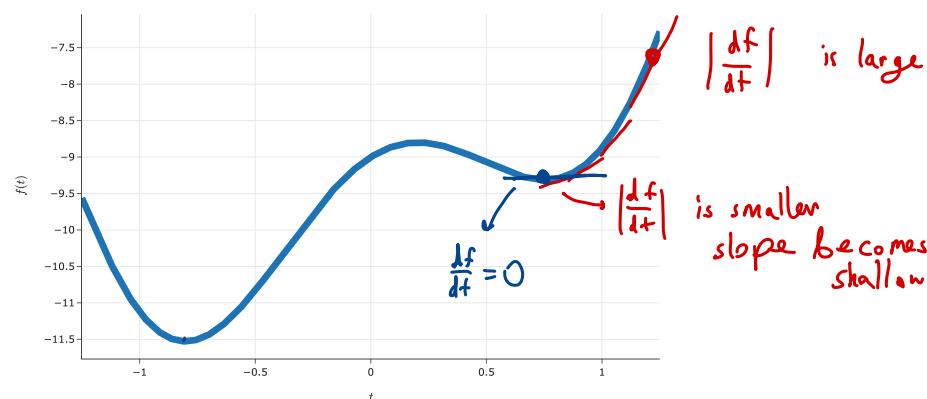
To minimize a **differentiable** function  $f$ :

- Pick a positive number,  $\alpha$ . This number is called the **learning rate**, or **step size**.
- Pick an **initial guess**,  $t_0$ .
- Then, repeatedly update your guess using the **update rule**:

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



- Repeat this process until **convergence** – that is, when  $t$  doesn't change much.



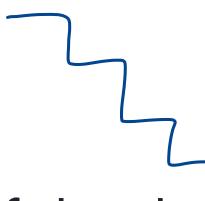
$$|t_n - t_{n-1}| < \varepsilon$$

# What is gradient descent?

- Gradient descent is a numerical method for finding the input to a function  $f$  that minimizes the function.
- Why is it called gradient descent?
  - The gradient is the extension of the derivative to functions of multiple variables.
  - We will see how to use gradient descent with multivariate functions next class.
- What is a numerical method?
  - A numerical method is a technique for approximating the solution to a mathematical problem, often by using the computer.

It is iterative

$t_0, t_1, t_2, t_3, \dots$



## Gradient descent

implementation of

$$\frac{df}{dh}$$

initialization  
( $t_0$ )

step size

convergence  
stopping criteria  
( $\epsilon$ )

```
def gradient_descent(derivative, h, alpha, tol=1e-12):  
    """Minimize using gradient descent."""
```

```
while True:
```

```
    h_next = h - alpha * derivative(h)
```

```
    if abs(h_next - h) < tol:
```

```
        break
```

```
    h = h_next
```

```
return h
```

$$h^* = \arg \min f(h)$$

stopping  
criteria

$$|h_{h+1} - h_h| < tol$$

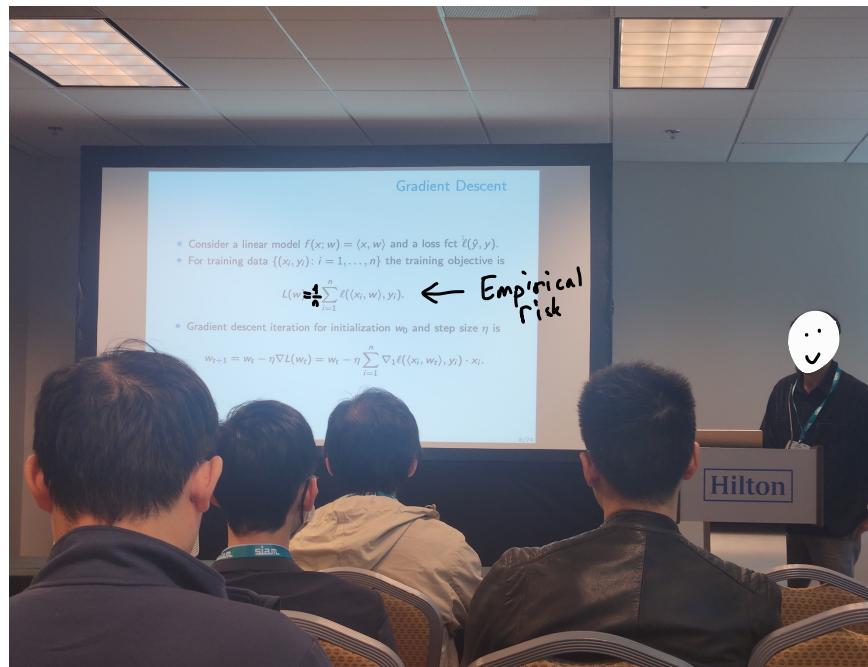
$$h_{n+1} = h_n - \alpha \cdot \frac{df}{dh}(h_n)$$

next point to go to  
we are at

See [this notebook](#) for a demo!

# Gradient descent and empirical risk minimization

- While gradient descent can minimize other kinds of differentiable functions, its most common use case is in **minimizing empirical risk**.
- Gradient descent is **widely used** in machine learning, to train models from linear regression to neural networks and transformers (includng ChatGPT)!



## Question 🤔

Answer at [q.dsc40a.com](http://q.dsc40a.com)

- For example, consider:
  - The constant model,  $H(x) = h$ .
  - The dataset  $-4, -2, 2, 4$ .
  - The initial guess  $h_0 = 4$  and the learning rate  $\alpha = \frac{1}{4}$ .

- Exercise: Find  $h_1$  and  $h_2$ .

$$h_1 = h_0 - \alpha \frac{dR_{sq}}{dh}(h_0)$$

$$h_2 = h_1 - \alpha \frac{dR_{sq}}{dh}(h_1)$$

$$l = (y_i - h)^2$$

$$R_{sq} = ?$$

$$\frac{dR_{sq}}{dh} = ?$$

typical  
exam  
question!

Common mistakes  
students make

\*  $- \rightarrow +$

\* forget  $\alpha$

\* plug in wrong  $h$  in  $\frac{dR_{sq}}{dh}(h)$

## Empirical Minimization with Gradient Descent

$$R_{\text{sq}} = \frac{1}{n} \sum_{i=1}^n (y_i - h)^2$$

$$\frac{dR_{\text{sq}}}{dh} = \frac{2}{n} \sum_{i=1}^n (h - y_i) = \frac{2}{n} nh - 2 \frac{1}{n} \sum_{i=1}^n y_i$$

$$= 2h - 2\bar{y}$$

- The dataset  $-4, -2, 2, 4$ .
- The initial guess  $h_0 = 4$  and the learning rate  $\alpha = \frac{1}{4}$ .

$$h_1 = h_0 - \alpha \frac{dR_{\text{sq}}}{dh}(h_0)$$

$$= 4 - \frac{1}{4} \cdot 8 = 2$$

$$\frac{dR_{\text{sq}}}{dh}(h_0) = \frac{2}{4} (4 - (-4) + 4 - (-2) + 4 - 2 + 4 - 4)$$

$$= \frac{1}{2} (8 + 6 + 2 + 0) = 8$$

$$h_2 = h_1 - \alpha \frac{dR_{\text{sq}}}{dh}(h_1)$$

$$= 2 - \frac{1}{4} \cdot 4 = 2 - 1 = 1$$

$$\frac{dR_{\text{sq}}}{dh}(h_1) = 2 \cdot 4 - 2 \cdot 0 = 8$$

$$\frac{dR_{\text{sq}}}{dh}(h_1) = 2 \cdot 2 - 2 \cdot 0 = 4$$

$4 \rightarrow 2 \rightarrow 1 \rightarrow \dots \rightarrow 0 = \bar{y}$

Midterm topics do not include:

\* center & spread (questions in practice site  
about mean absolute deviation)

\* gradient descent

HW4 solution will be released on Sunday

## Lingering questions

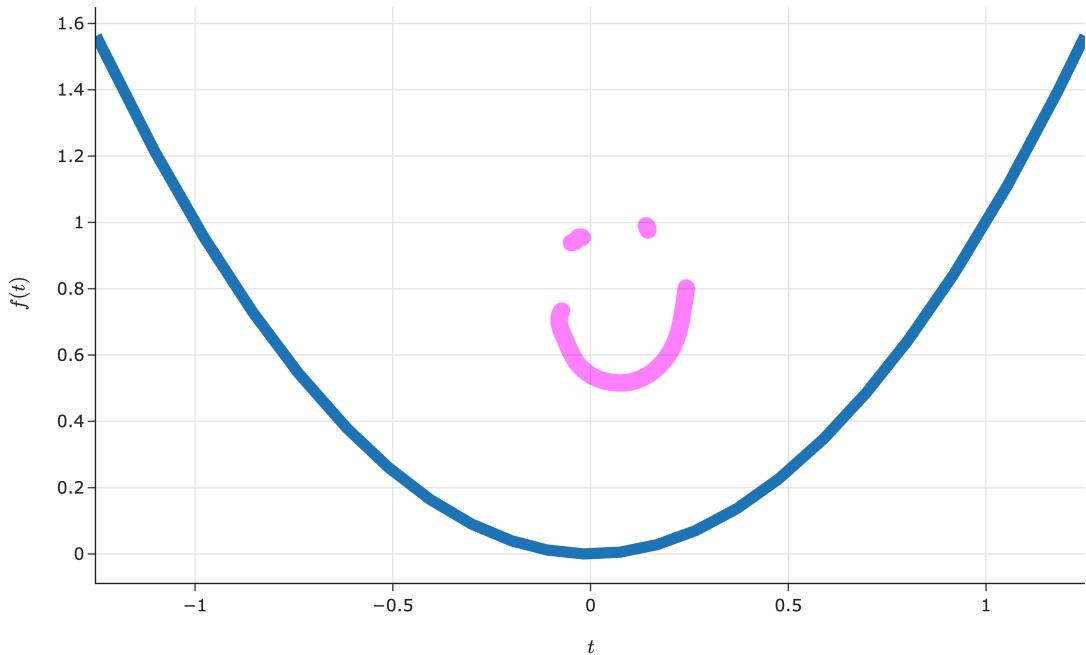
Now, we'll explore the following ideas:

- When is gradient descent *guaranteed* to converge to a global minimum?
  - What kinds of functions work well with gradient descent?
- How do I choose a step size?
- How do I use gradient descent to minimize functions of multiple variables, e.g.:

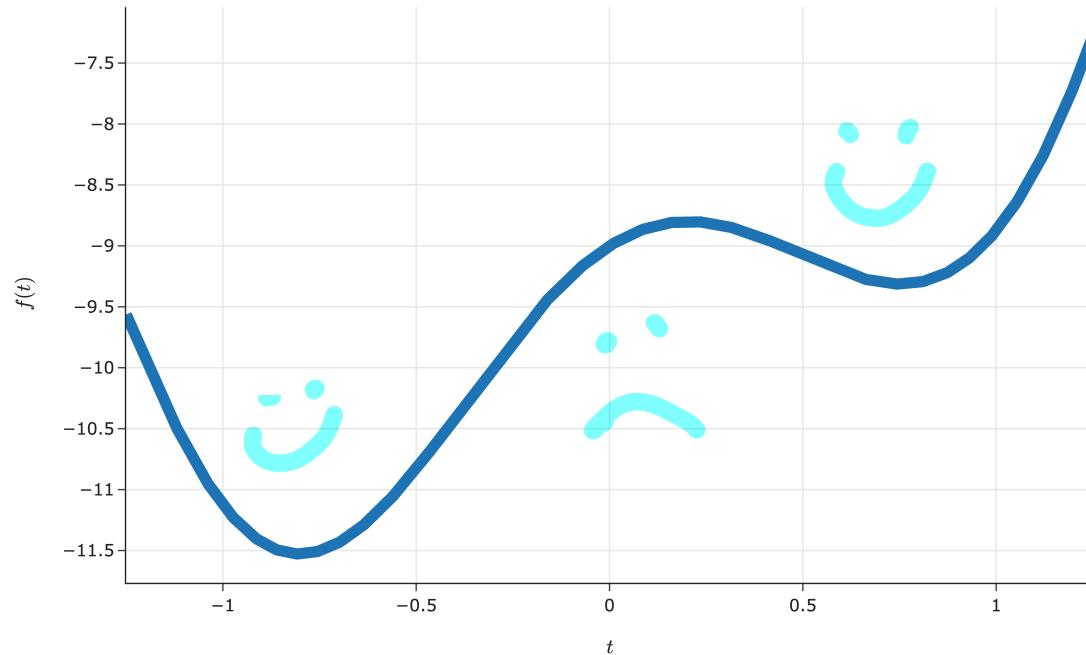
$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

# When is gradient descent guaranteed to work?

# Convex functions



A convex function



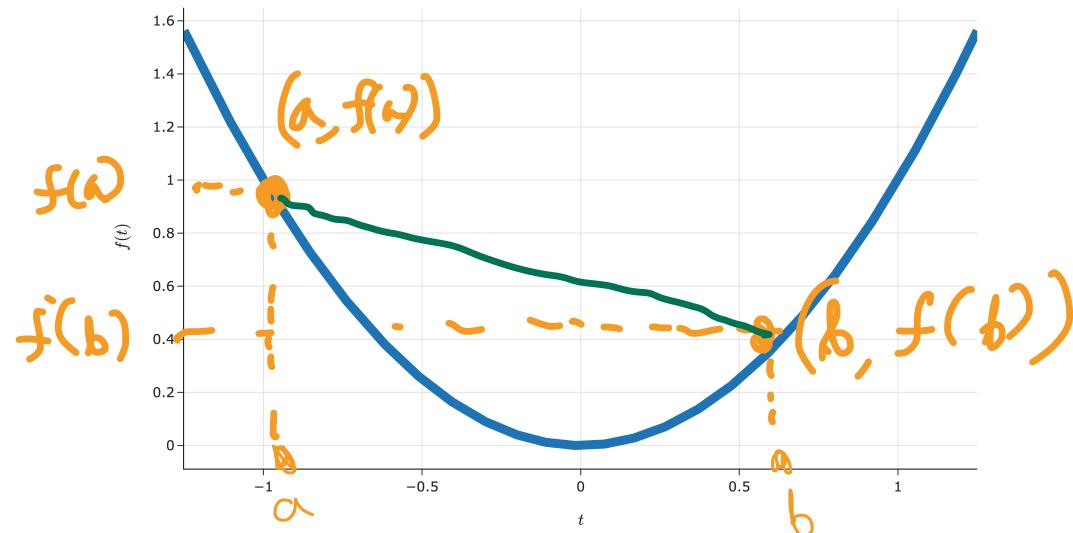
A non-convex function

# Convexity

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

$(a, f(a))$  and  $(b, f(b))$

does not go below the plot of  $f$ .



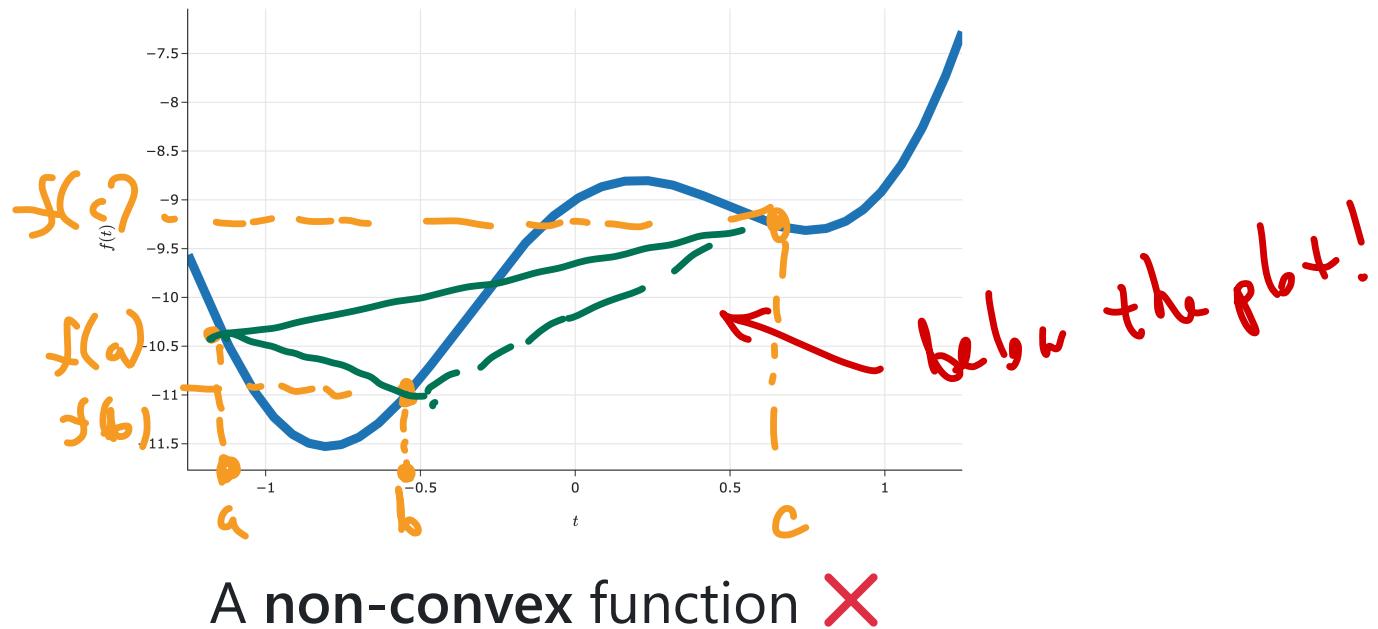
A convex function

# Convexity

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

$(a, f(a))$  and  $(b, f(b))$

does not go below the plot of  $f$ .



## Formal definition of convexity

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

plug in  $t=0$       plug in  $t=1$   
 $f(a)$                    $f(b)$

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

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

Ex:  $t = \frac{1}{2}$

$$\frac{1}{2}f(a) + \frac{1}{2}f(b) \geq f\left(\frac{1}{2}a + \frac{1}{2}b\right)$$

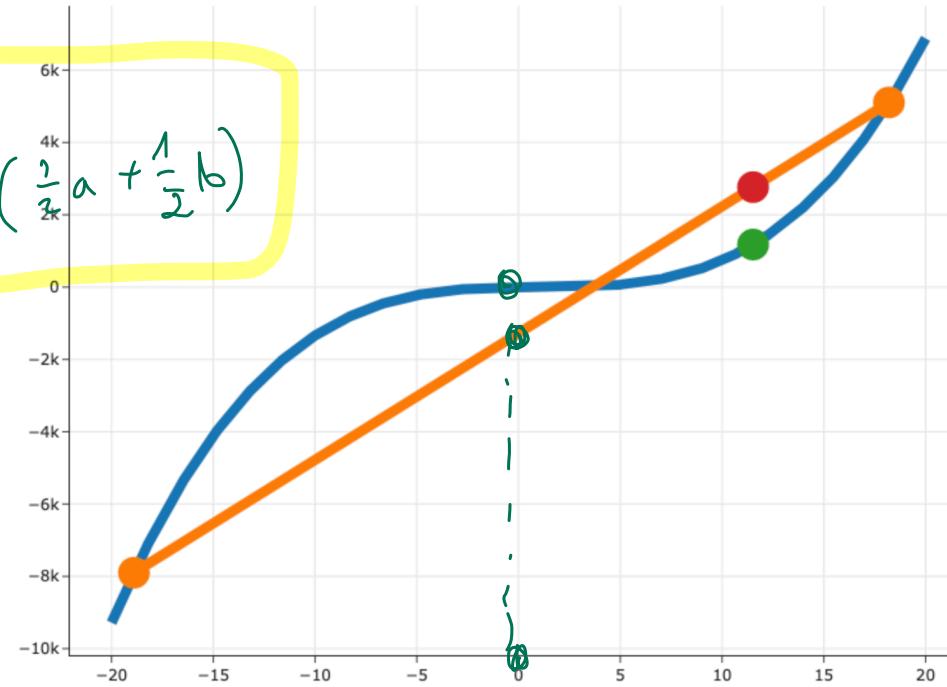
Function between  
 $x=a$  and  $x=b$

- A function is nonconvex if it is not convex.
- This is a formal way of restating the definition from the previous slide.

If  $0 \leq t \leq 1$  what is

$$50 + 30t$$

$$(1-t)50 + 90t \quad 0 \leq t \leq 1$$



$$\frac{1}{2}a + \frac{1}{2}b$$

line  $\geq$  function  
 for  $0 \leq t \leq 1$

$\iff$  Convex

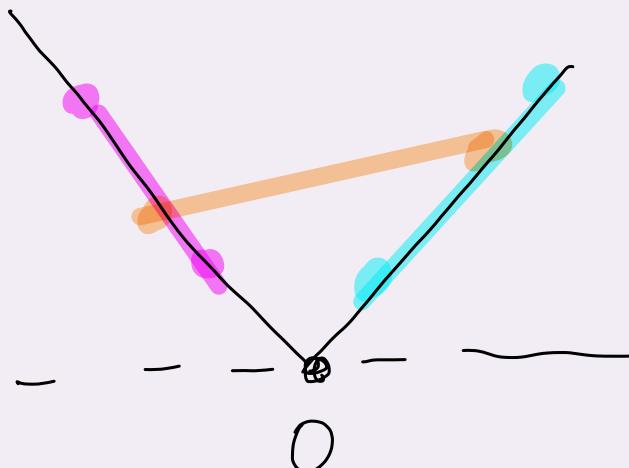
$$f(x) = |x|$$

## Question 🤔

Answer at [q.dsc40a.com](http://q.dsc40a.com)

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

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



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

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

$$(1-t)f(a) + t f(b) \geq f((1-t)a + t(b)) \quad \text{for all } 0 \leq t \leq 1$$

$$(1-t)|a| + t|b| \geq |(1-t)a + t(b)|$$

$$|(1-t)a + t(b)| \leq |(1-t)a| + |t(b)| \leq (1-t)|a| + t|b|$$

the segment

always non-negative  
for  $0 \leq t \leq 1$

function

## Question 🤔

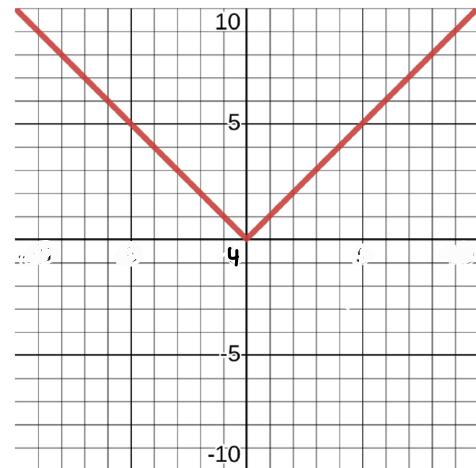
Answer at [q.dsc40a.com](http://q.dsc40a.com)

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.

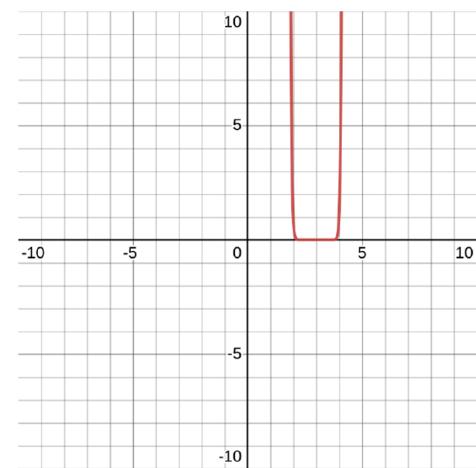
# Convex vs. concave

Convex



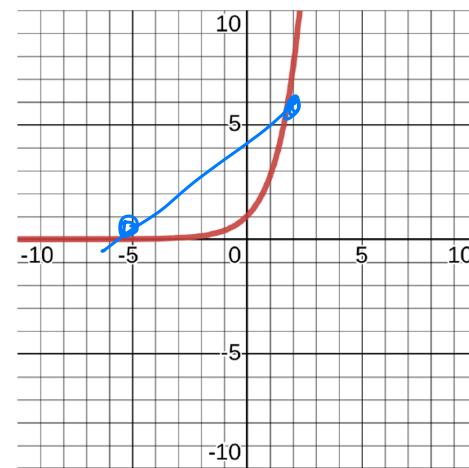
$$f(x) = |x - 4|$$

Convex



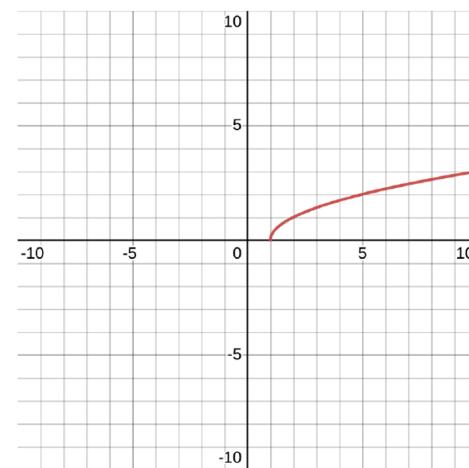
$$f(x) = (x - 3)^2$$

Convex



$$f(x) = e^x$$

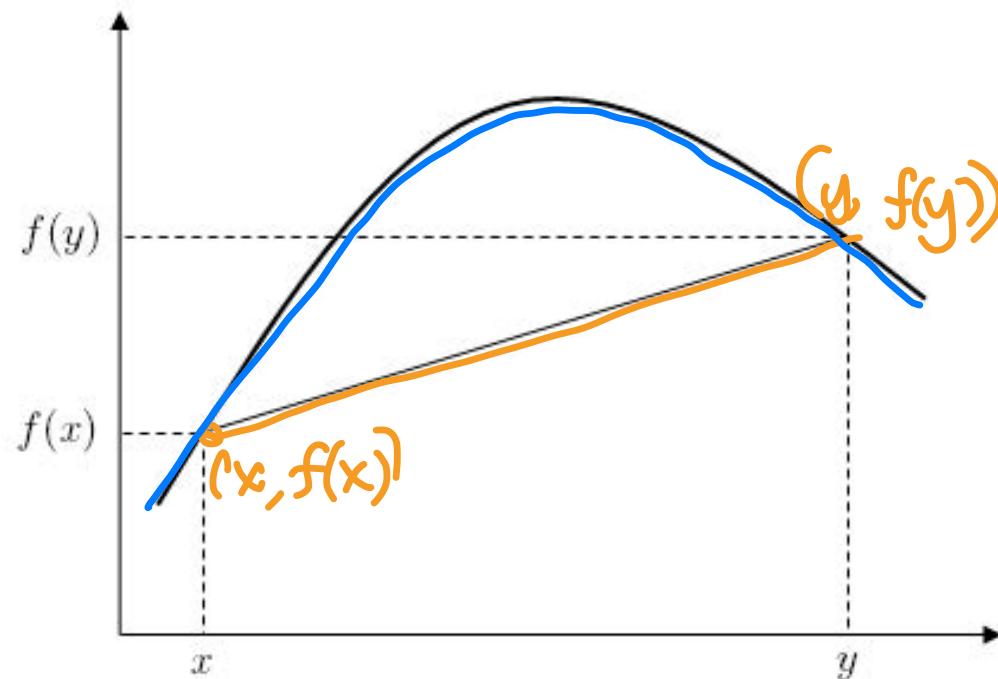
Concav



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

# 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:



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

↙ for all t

- concave if and only if:

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

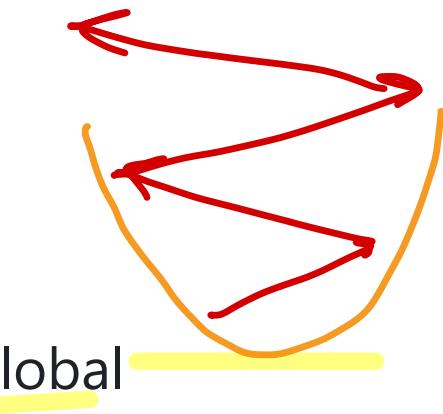
- Example:  $f(x) = x^4$  is convex.

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

$$f''(x) = 12x^2 \geq 0 \quad \forall x \Rightarrow \text{convex}$$

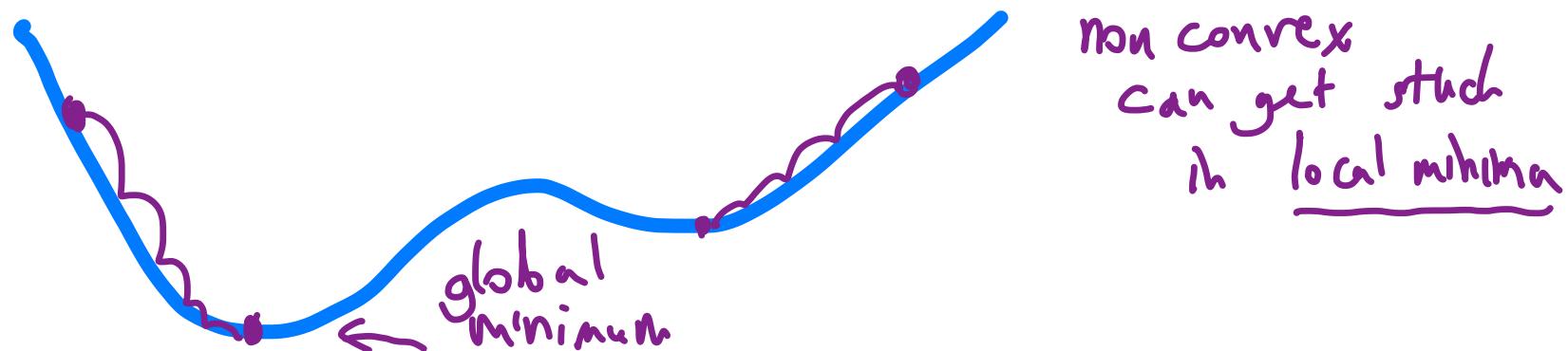
# Why does convexity matter?

- 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.).



# 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
$$f(t) = 5t^4 - t^3 - 5t^2 + 2t - 9.$$



## 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

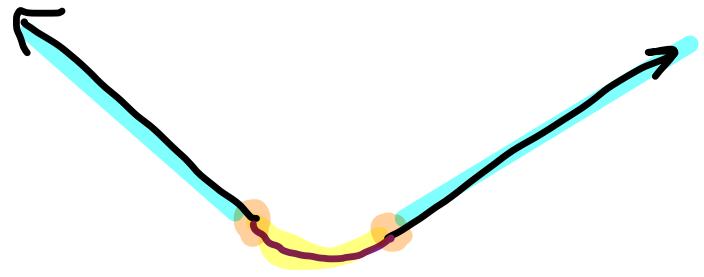
## Example: Huber loss and the constant model

- 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,

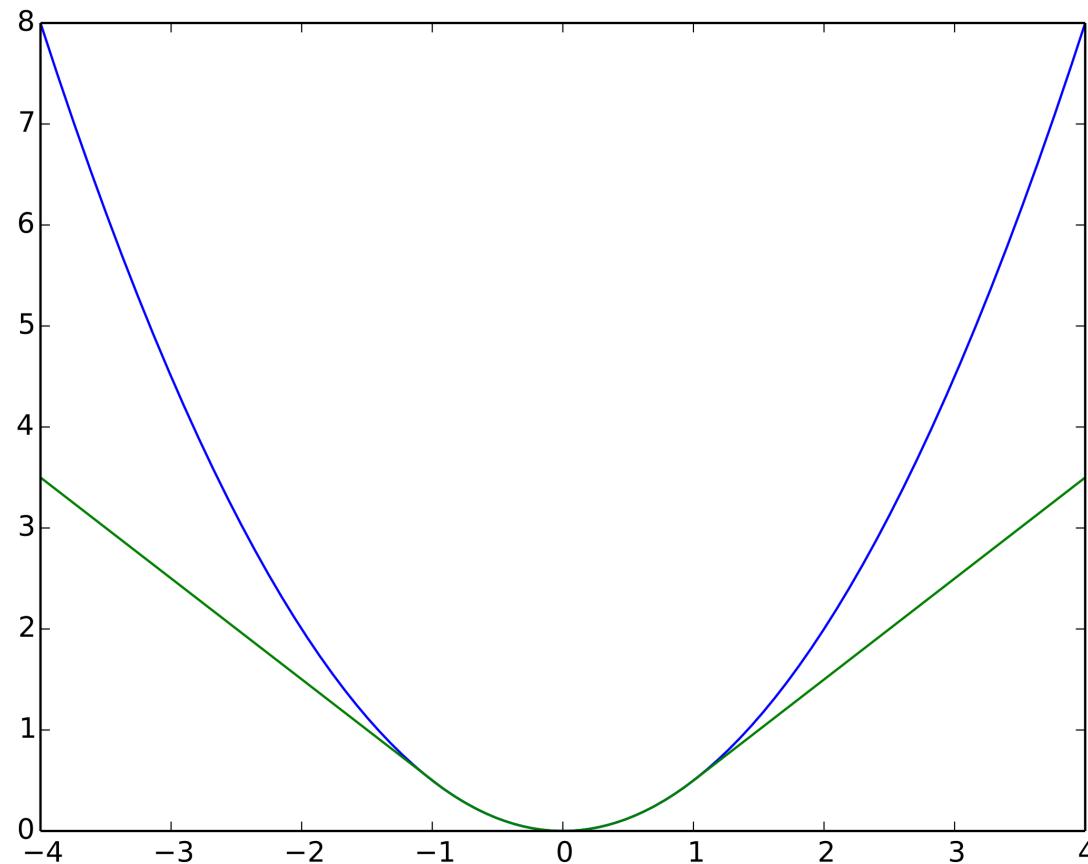
$$L_{\text{abs}}(y_i, H(x_i)) = |y_i - H(x_i)|.$$

pro: robust to outliers

con: not differentiable, harder to minimize

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

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



**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,  $H(x) = h$ :

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

$$\begin{aligned} y_i &= h \\ \frac{\partial L}{\partial h}(h) &= \begin{cases} 0 & \text{if } y_i = h \\ -\delta & \text{otherwise} \end{cases} = 0 \end{aligned}$$

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

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

- 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!

Let's try this out in practice! Follow along in [this notebook](#).

## Minimizing functions of multiple variables

- Consider the function:

$$f(x_1, x_2) = (x_1 - 2)^2 + 2x_1 + (x_2 - 3)^2$$

- It has two **partial derivatives**:  $\frac{\partial f}{\partial x_1}$  and  $\frac{\partial f}{\partial x_2}$ .

## The gradient vector

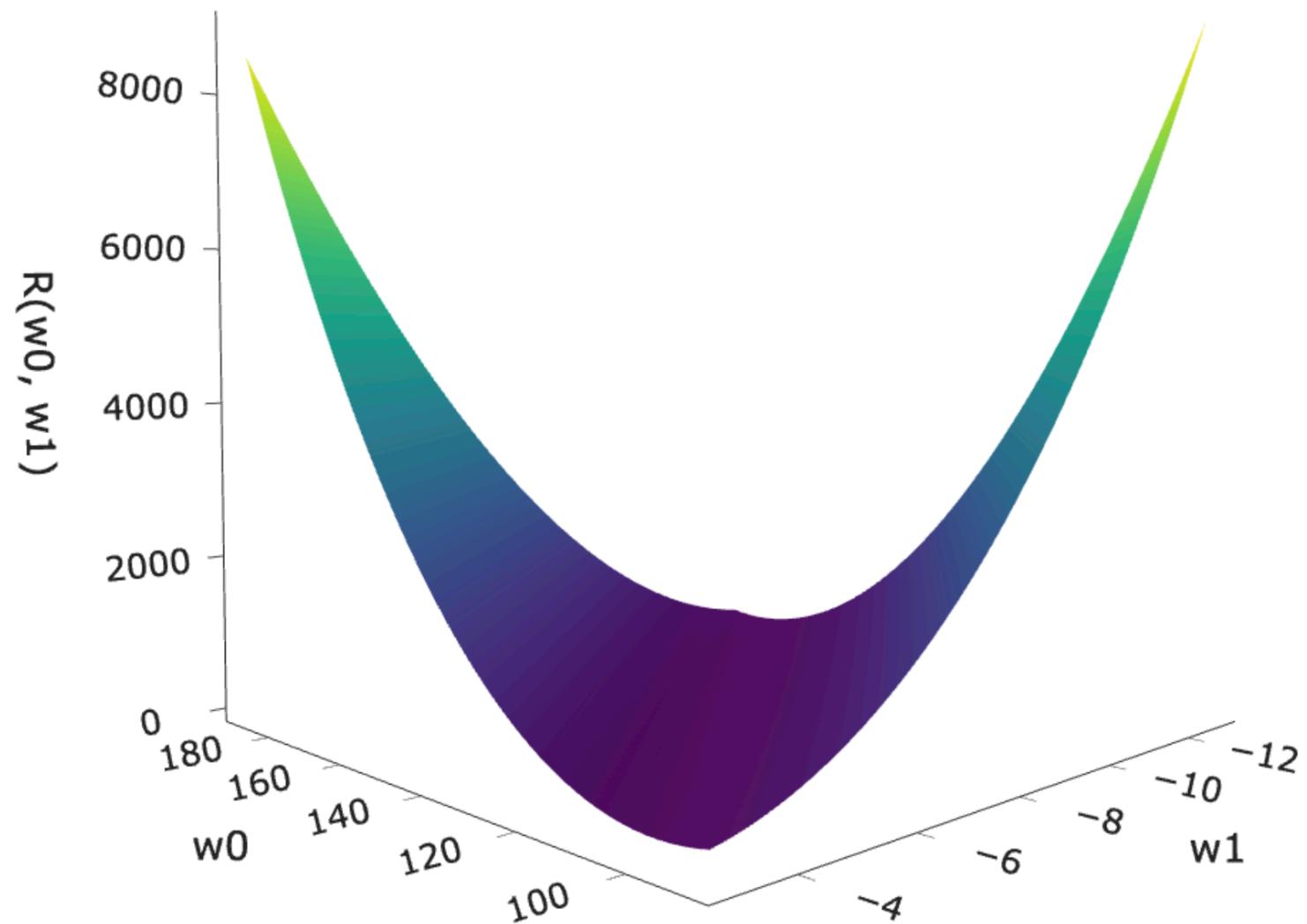
- If  $f(\vec{x})$  is a function of multiple variables, then its **gradient**,  $\nabla f(\vec{x})$ , is a vector containing its partial derivatives.
- Example:

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

$$\nabla f(\vec{x}) = \begin{bmatrix} 2x_1 - 2 \\ 2x_2 - 6 \end{bmatrix}$$

- Example:

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



## Gradient descent for functions of multiple variables

- Example:

$$f(x_1, x_2) = (x_1 - 2)^2 + 2x_1 + (x_2 - 3)^2$$

$$\nabla f(\vec{x}) = \begin{bmatrix} 2x_1 - 2 \\ 2x_2 - 6 \end{bmatrix}$$

- The minimizer of  $f$  is a vector,  $\vec{x}^* = \begin{bmatrix} x_1^* \\ x_2^* \end{bmatrix}$ .
- We start with an initial guess,  $\vec{x}^{(0)}$ , and step size  $\alpha$ , and update our guesses using:

$$\vec{x}^{(i+1)} = \vec{x}^{(i)} - \alpha \nabla f(\vec{x}^{(i)})$$

## Exercise

$$f(x_1, x_2) = (x_1 - 2)^2 + 2x_1 + (x_2 - 3)^2$$

$$\nabla f(\vec{x}) = \begin{bmatrix} 2x_1 - 2 \\ 2x_2 - 6 \end{bmatrix}$$

$$\vec{x}^{(i+1)} = \vec{x}^{(i)} - \alpha \nabla f(\vec{x}^{(i)})$$

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** iterations of gradient descent. What is  $\vec{x}^{(2)}$ ?



## Example: Gradient descent for simple linear regression

- To find optimal model parameters for the model  $H(x) = w_0 + w_1x$  and squared loss, we minimized empirical risk:

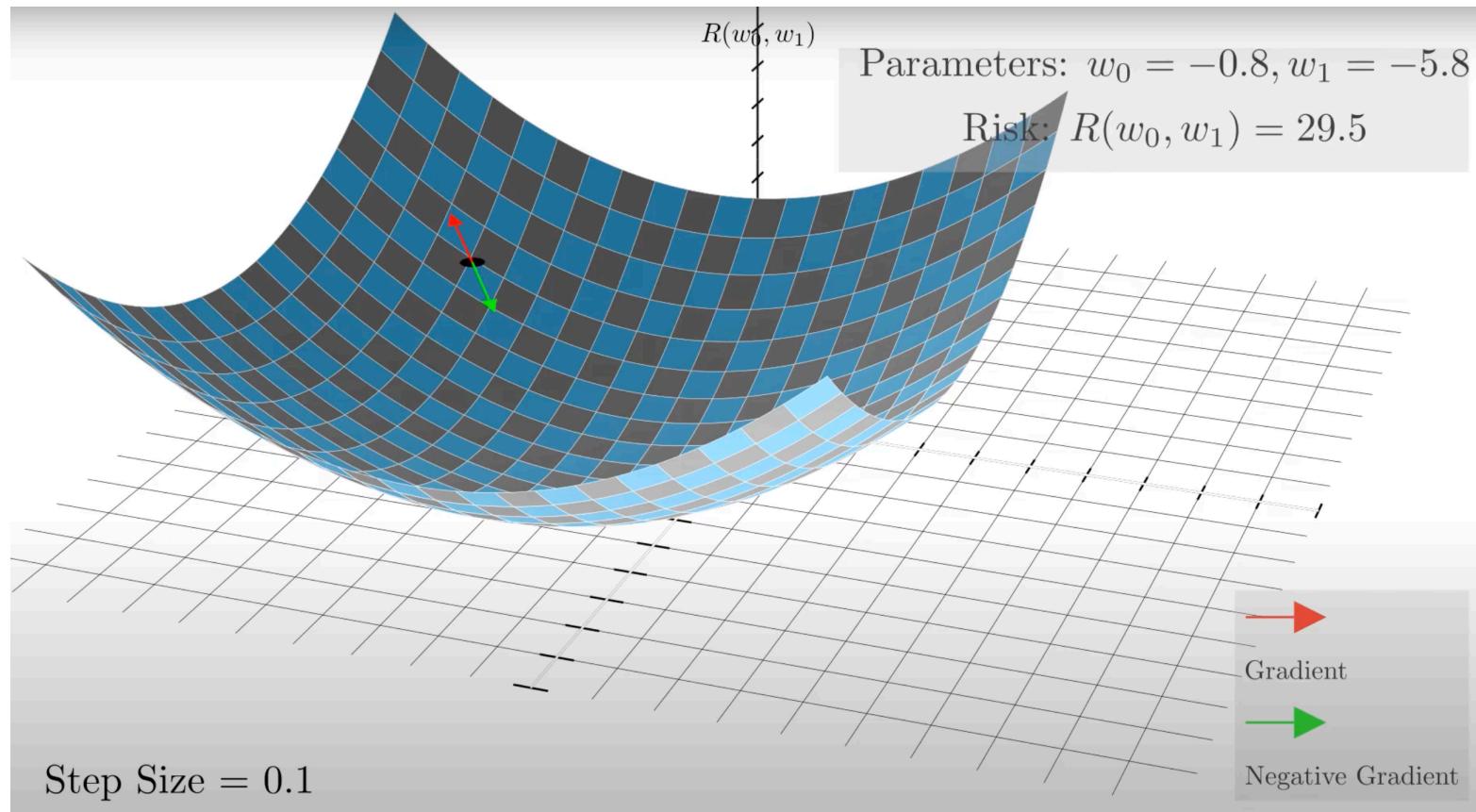
$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

- This is a function of multiple variables, and is differentiable, so it has a gradient!

$$\nabla R(\vec{w}) = \begin{bmatrix} -\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) \\ -\frac{2}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i)) x_i \end{bmatrix}$$

- Key idea:** To find  $w_0^*$  and  $w_1^*$ , we *could* use gradient descent!

# Gradient descent for simple linear regression, visualized



Let's watch [this animation](#) that Jack made.

## What's next?

- In Homework 5, you'll see a few questions involving today's material.
- After the midterm, we'll start talking about probability.