

Lecture 11

Gradient Descent, Continued

DSC 40A, Summer 2024

Announcements

- The Midterm Exam is **tomorrow!**
- Some time for review in Discussion today, and Owen's OH 4-~~5~~⁶p in HDSI 155.

Agenda

- Recap: Gradient descent.
- Convexity.
- More examples.
 - Huber loss.
 - Gradient descent with multiple variables.

Question 🤔

Answer at q.dsc40a.com

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

If the direct link doesn't work, click the " Lecture Questions" link in the top right corner of dsc40a.com.

Overview: Gradient descent

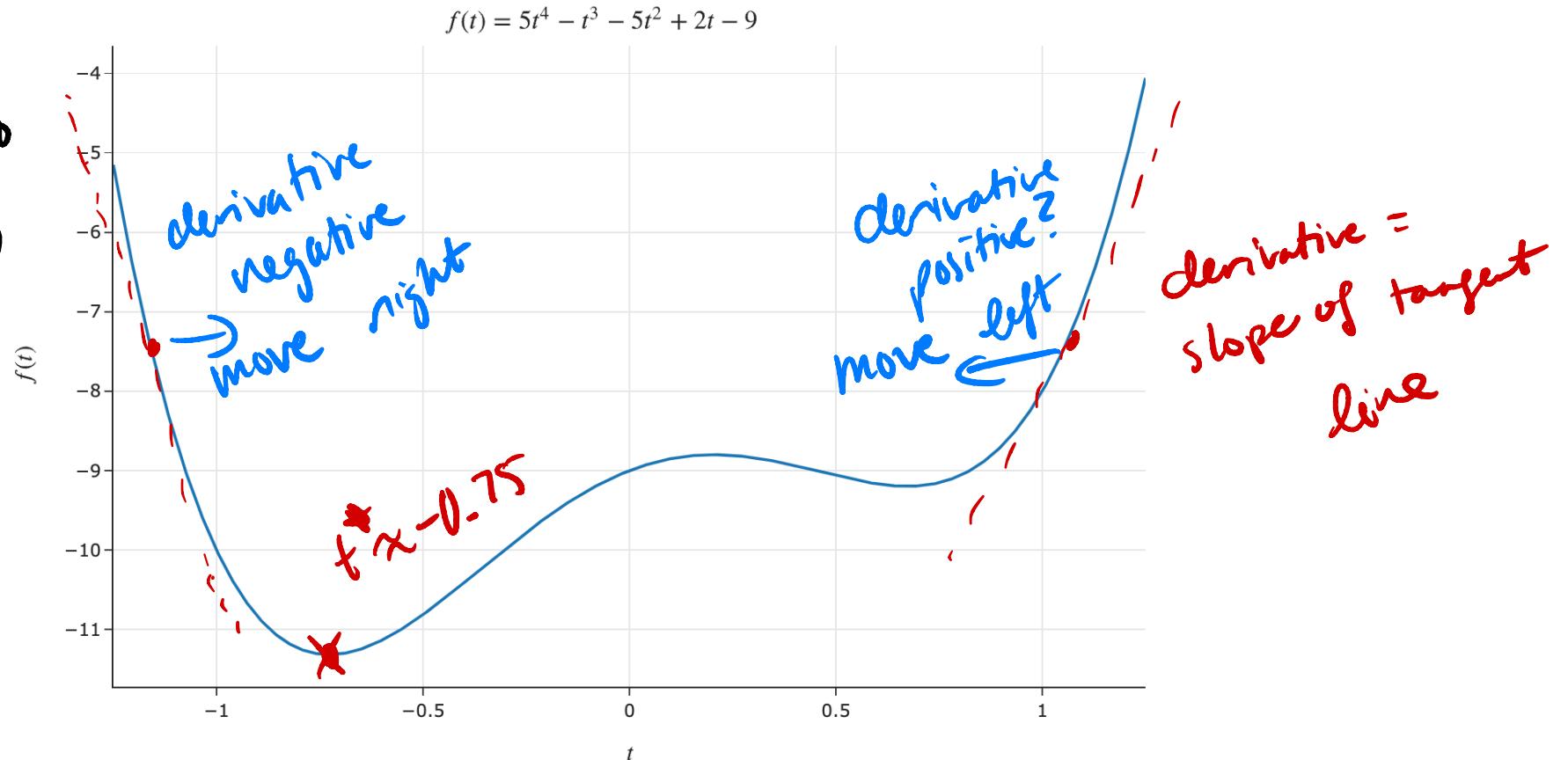
derivative exists everywhere

What's the point?



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

Can find $\frac{df}{dt}$
but solve for $\frac{df}{dt}(t) = 0$
maybe hard to



Gradient descent

To minimize a **differentiable** function f :

- Pick a positive number, α . 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**:

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

learning rate / step size

negative, because we move opposite the derivative

- Repeat this process until **convergence** – that is, when t doesn't change much.
- This procedure is called **gradient descent**.

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 ~~this~~ ^{this} 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.
- Gradient descent is **widely used** in machine learning, to train models from linear regression to neural networks and transformers (including ChatGPT)!

descending to the bottom

See dsc40a.com/resources/lectures/lec10 for animated examples of gradient descent, and see [this notebook](#) for the associated code!

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**.
- For example, consider:
 - The constant model, $H(x) = h$.
 - Squared loss.
 - 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 .

$$\text{data} = -4, -2, 2, 4$$

$$h_0 = 4, \alpha = \frac{1}{4}$$

$$h_0 = 4$$
$$h_1 = h_0 - \alpha \frac{d R_{SG}(h)}{dh} \quad \boxed{2h_0}$$

$$= 4 - \frac{1}{4} (2 \cdot 4)$$
$$= 4 - 2 = \boxed{2 = h_1}$$

$$h_2 = h_1 - \alpha (2 h_1)$$
$$= 2 - \frac{1}{4} \cdot (2 \cdot 2)$$
$$= 2 - 1 = \boxed{1 = h_2}$$

$$R_{SG}(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h)^2$$
$$\frac{d R_{SG}(h)}{dh} = \frac{1}{n} \sum_{i=1}^n 2(y_i - h)(-1)$$
$$= -\frac{2}{n} \sum_{i=1}^n (y_i - h)$$
$$-4 + (-2) + 2 + 4 = 0$$
$$= -\frac{2}{n} \left[\sum_{i=1}^n (y_i - h) \right]$$

$$\frac{d R_{SG}(h)}{h} = -\frac{2}{4} (-4h) = 2h$$

$$\boxed{\begin{array}{l} h_0 = 4 \\ h_1 = 2 \\ h_2 = 1 \end{array}}$$

Lingering questions

Now, we'll explore the following ideas:

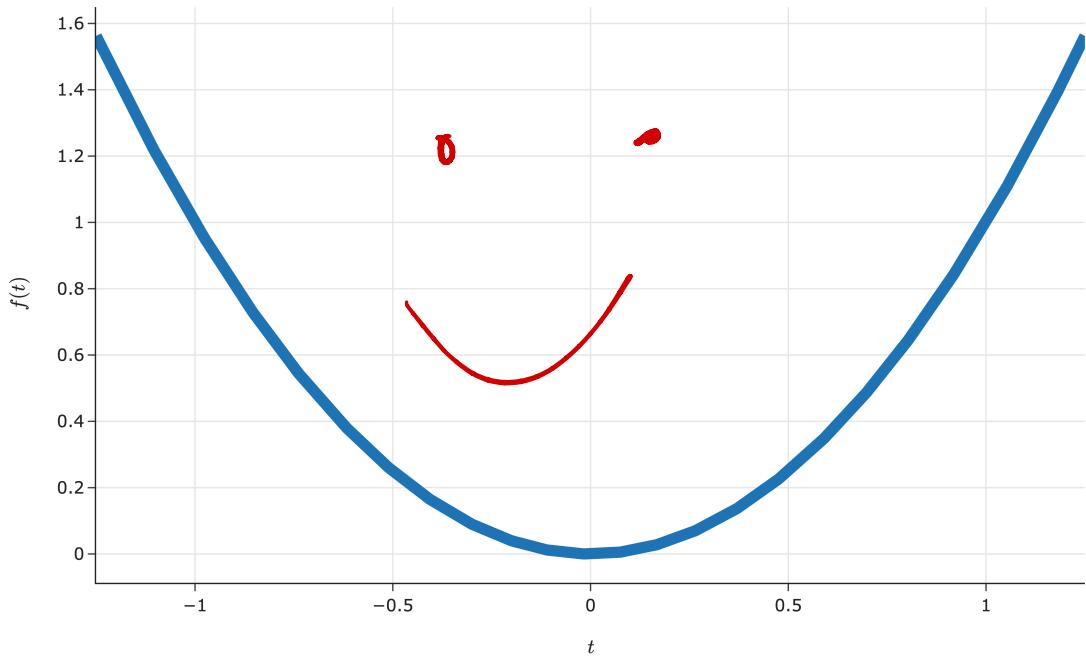
- When is gradient descent *guaranteed* to converge to a global minimum?
 - What kinds of functions work well with gradient descent?
- How do I choose a step size?
- How do I use gradient descent to minimize functions of multiple variables, e.g.:

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

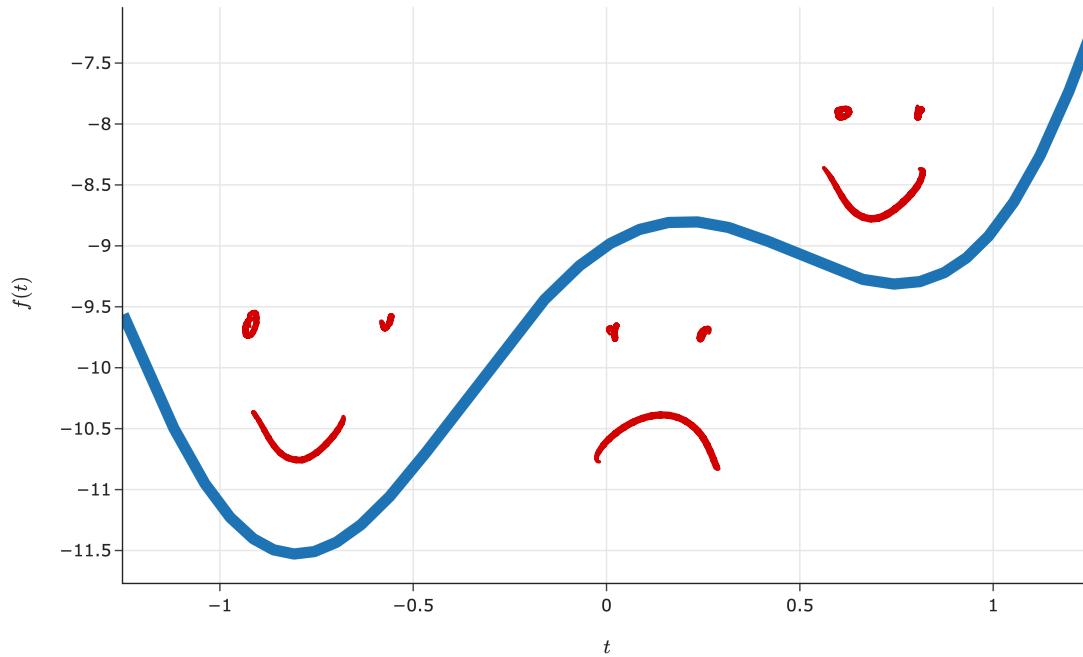
When is gradient descent guaranteed to work?

Convex functions

" bowl-shaped"



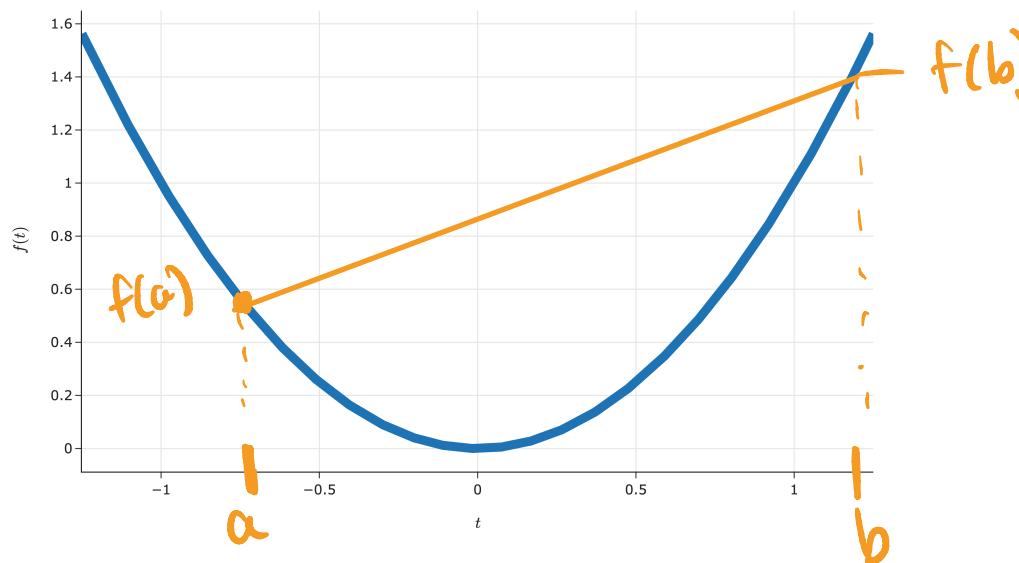
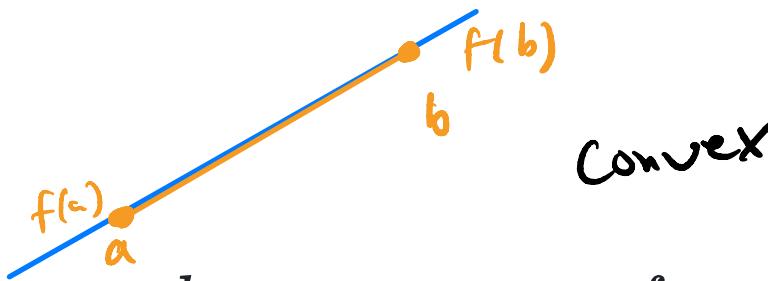
A convex function



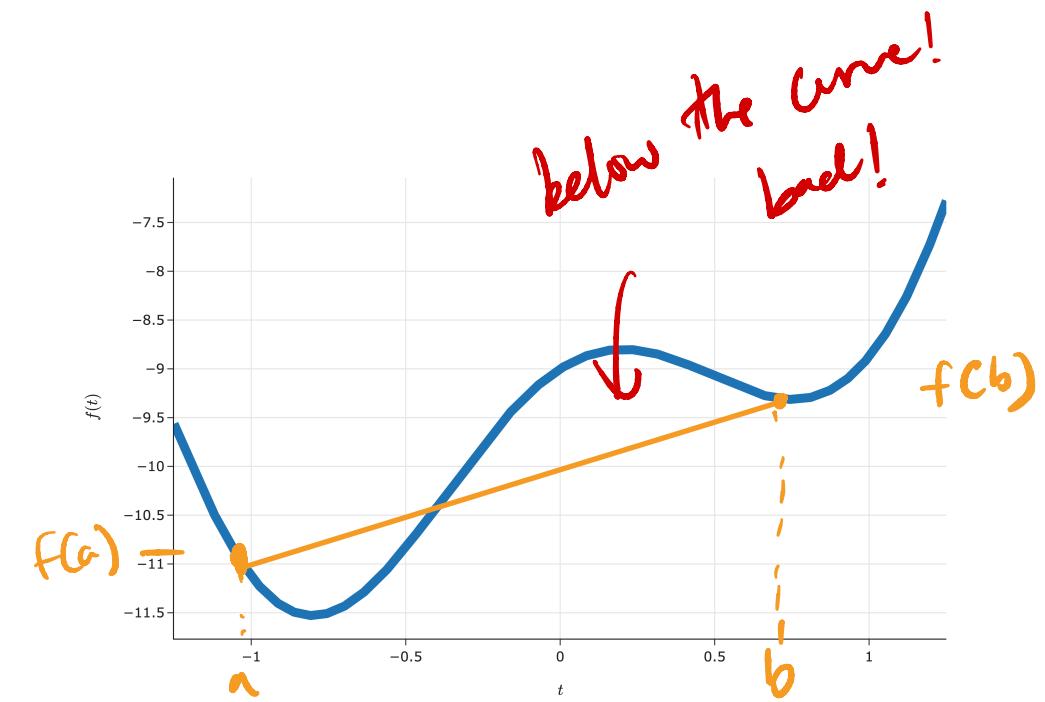
A non-convex function

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



A non-convex function

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 $t \in [0, 1]$:

Ex

plug in some linear function $y=mx+b$

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

function between $x=a$, and $x=b$

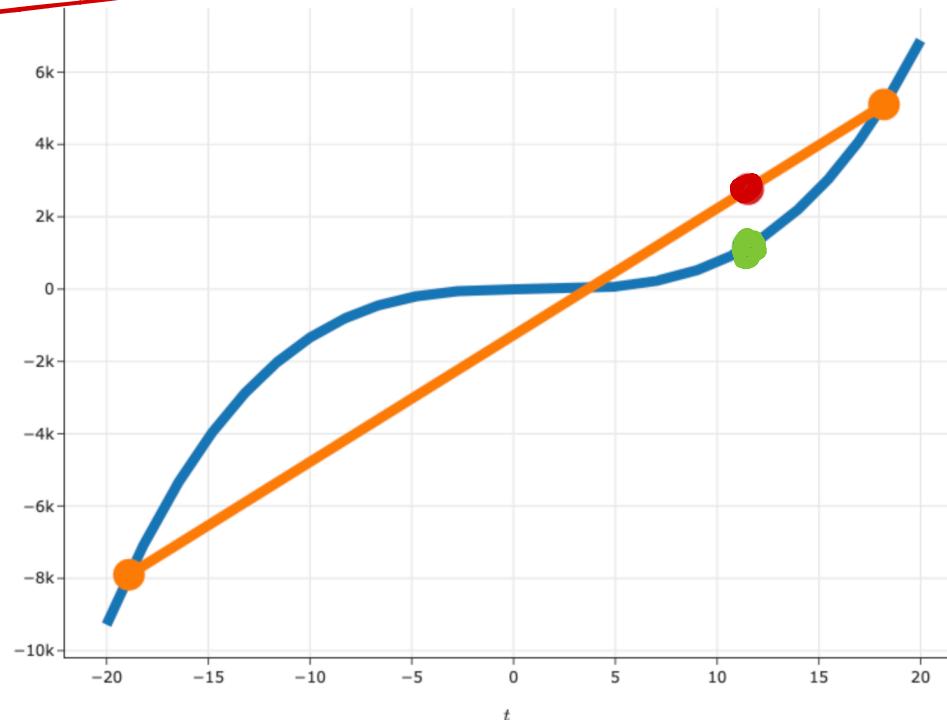
line between $f(a)$, and $f(b)$

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

need: line \geq function

function is convex

aside if $0 \leq t \leq 1$, what is
 $50 + 30t$
 \downarrow
 $= (1-t)50 + t(80)$



"convex combination" of a bb

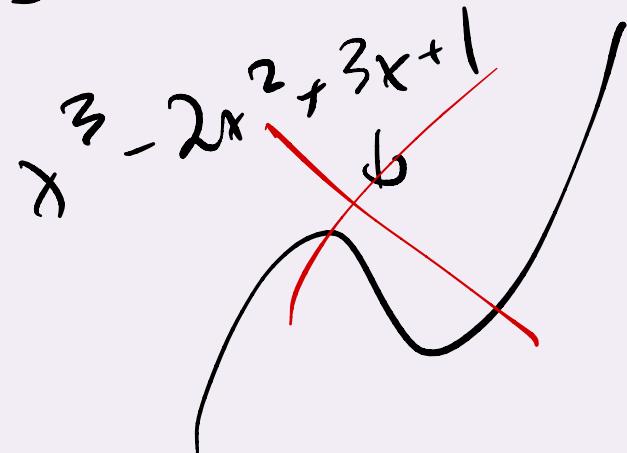
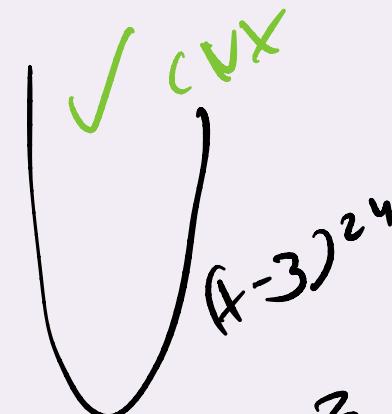
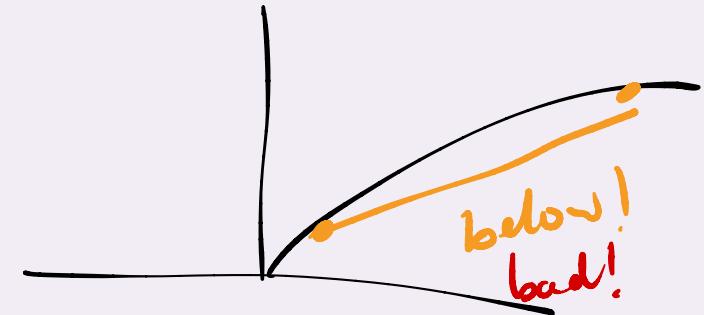
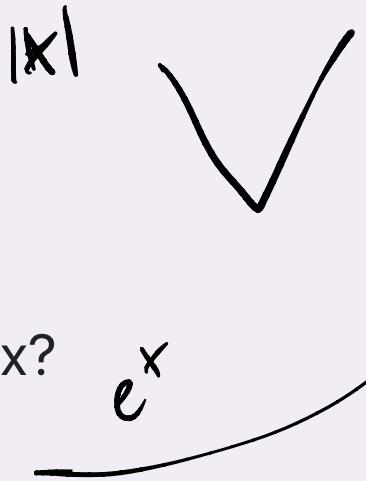
Question 🤔

Answer at q.dsc40a.com

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

- A. $f(x) = |x|$.
- B. $f(x) = e^x$.
- C. $f(x) = \sqrt{x - 1}$.
- D. $f(x) = (x - 3)^{24}$. → even power, like a parabola
- E. More than one of the above are non-convex.

$$\left\{ ((x-3)^2)^{12} \right\}$$



$\exists t$ "there exists a t "
 ↗
 single

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:

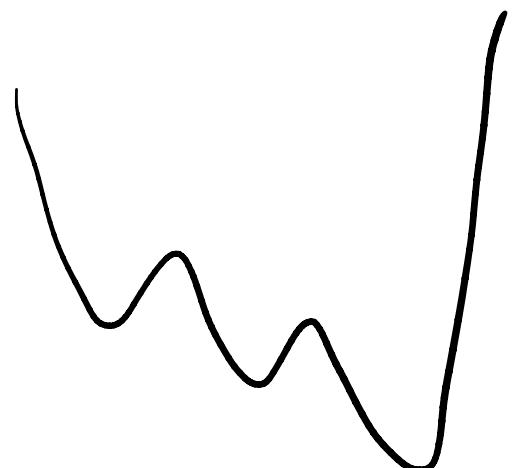
$$\frac{d^2 f}{dt^2}(t) \geq 0, \quad \forall t \quad \text{"for all } t\text{"}$$

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

x^2 also convex

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

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



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

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

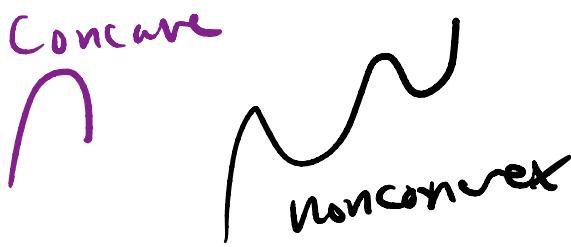
$$x^4 - 3x^3 + 2x^2 - 5x$$

$$4x^3 - 9x^2 + 4x - 5$$

$$12x^2 - 18x + 4 \stackrel{?}{\geq} 0 \quad \forall x$$

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.).

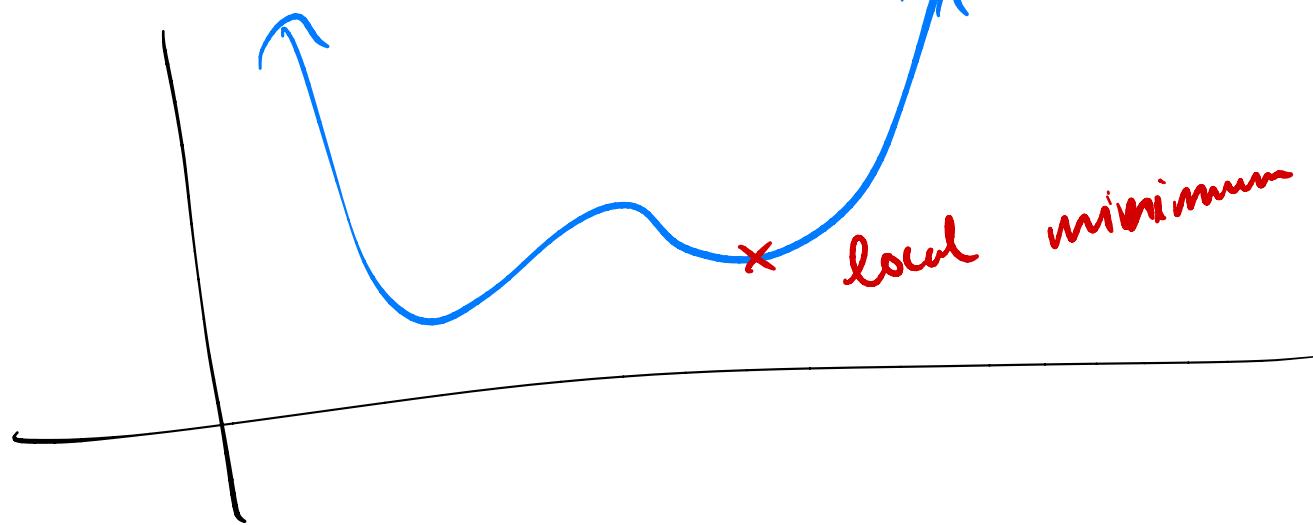


Nonconvex functions and gradient descent

do not say concave

- 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 α is a constant.

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

- **Remember:** α is the "step size", but the amount that our guess for t changes is $\alpha \frac{df}{dt}(t_i)$, not just α .
- In future courses, you'll learn about "decaying" step sizes, where the value of α 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: easy to minimize, differentiable

Con: sensitive to outliers

- Then, we learned about absolute loss,

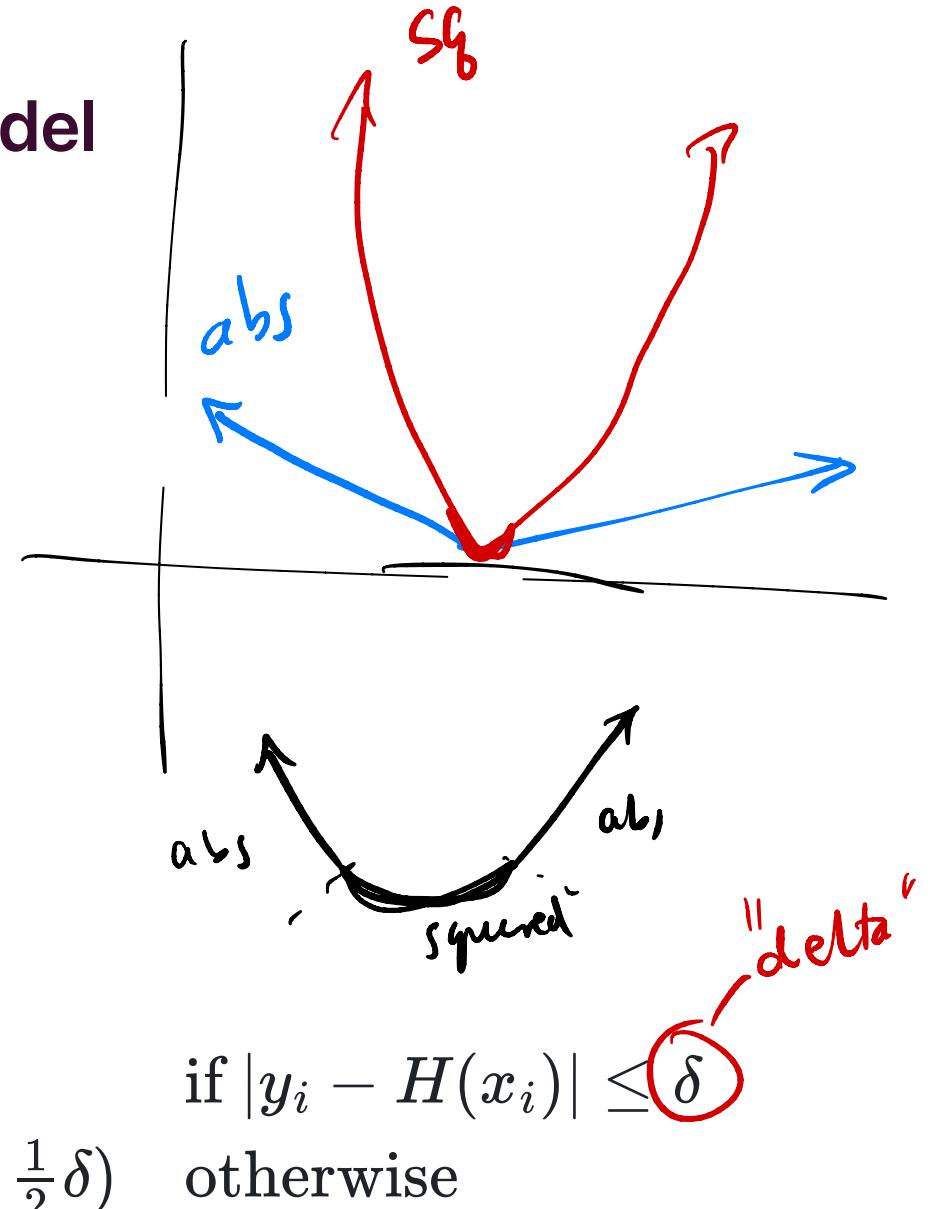
$$L_{\text{abs}}(y_i, H(x_i)) = |y_i - H(x_i)|.$$

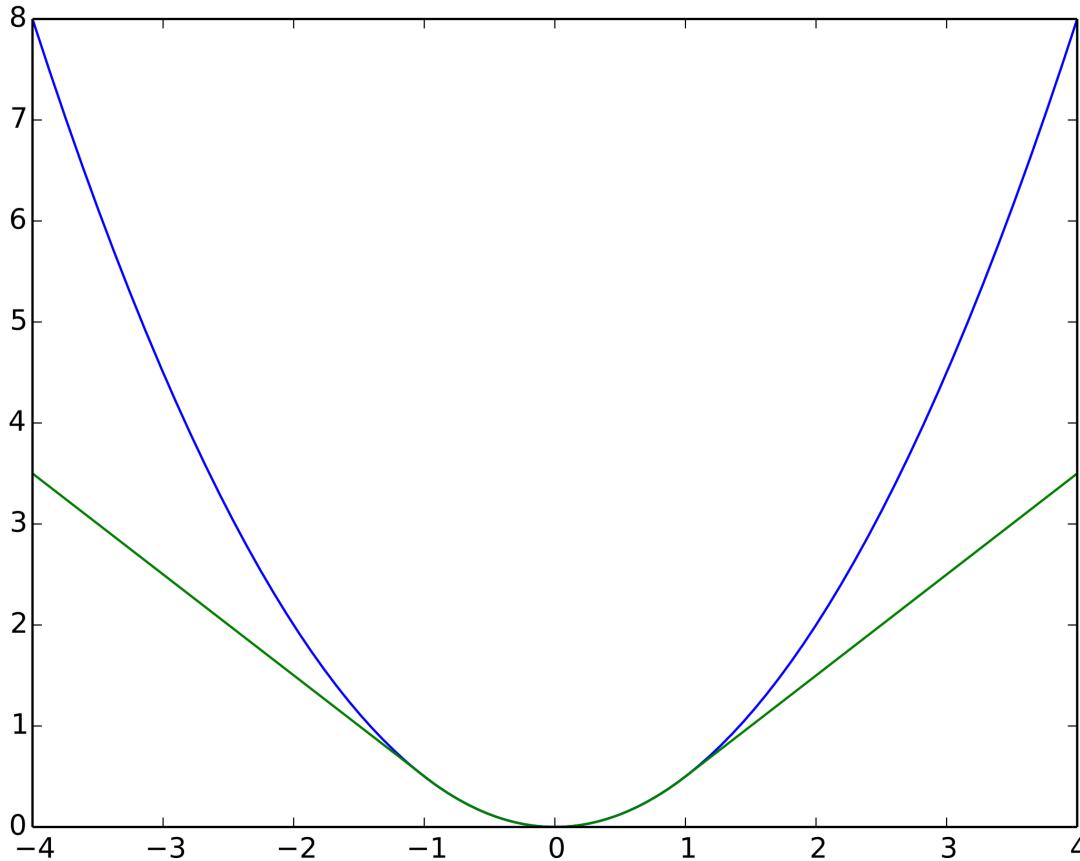
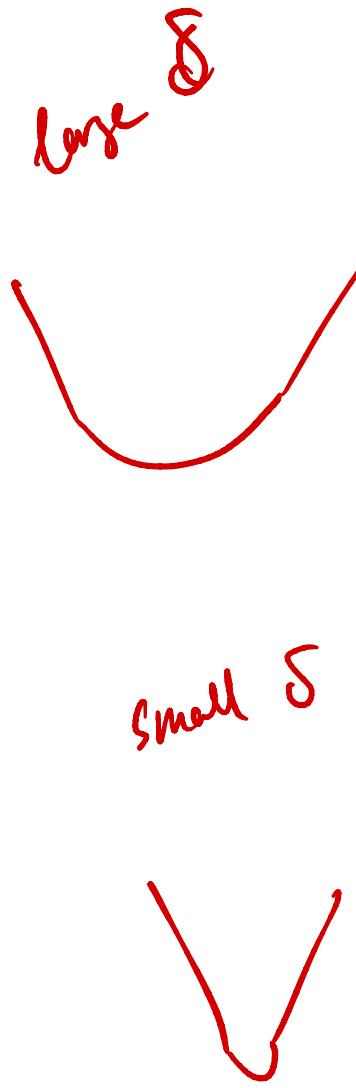
Pro: robust to outliers

Con: not differentiable

- 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, $H(x) = h$:

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

- 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 $\frac{\partial f}{\partial x_2}$.

$$\frac{\partial f}{\partial x_1} = 2x_1 - 2$$

$$\frac{\partial f}{\partial x_2} = -2(x_2 - 3)$$

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

"nabla" ∇

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:

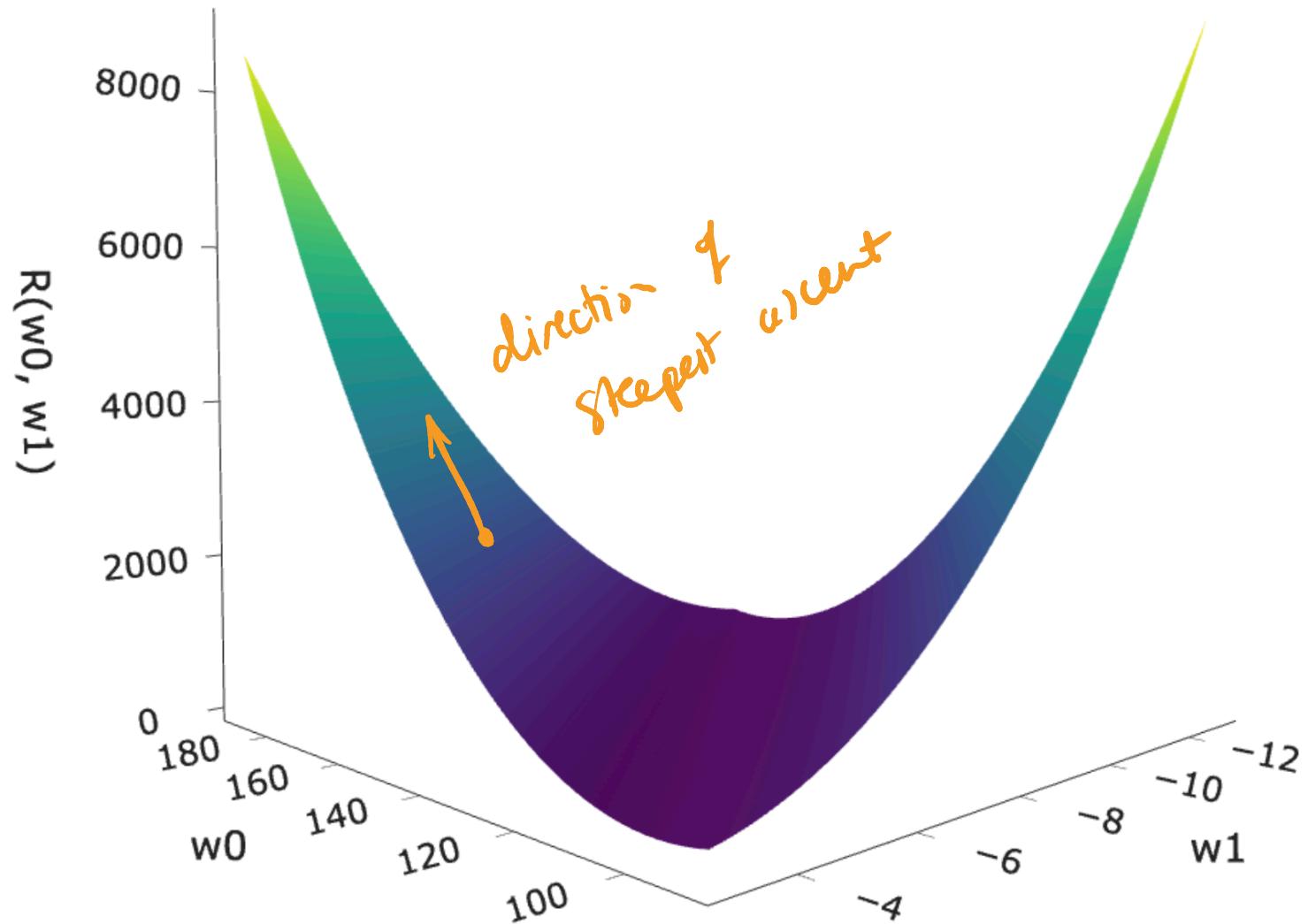
$$f(\vec{x}) = (x_1 - 2)^2 + 2x_1 - (x_2 - 3)^2$$

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

- Example:

$$f(\vec{x}) = \vec{x}^T \vec{x} = x_1^2 + x_2^2 + \dots + x_n^2$$
$$\Rightarrow \nabla f(\vec{x}) = \begin{bmatrix} 2x_1 \\ 2x_2 \\ \vdots \\ 2x_n \end{bmatrix} = 2\vec{x}$$

$\frac{\delta f}{\delta x_i} = 2x_i$



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 α , 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 $\vec{x}^{(2)}$?

$$\vec{x}^{(1)} = \vec{x}^{(0)} - \alpha \nabla f(\vec{x}^{(0)}) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} -2 \\ 6 \end{bmatrix} = \begin{bmatrix} 2/3 \\ -2 \end{bmatrix}$$

$$\vec{x}^{(2)} = \vec{x}^{(1)} - \alpha \nabla f(\vec{x}^{(1)}) = \begin{bmatrix} 2/3 \\ -2 \end{bmatrix} - \frac{1}{3} \begin{bmatrix} \frac{4}{3} - 2 \\ -(-4 - 6) \end{bmatrix}$$

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:

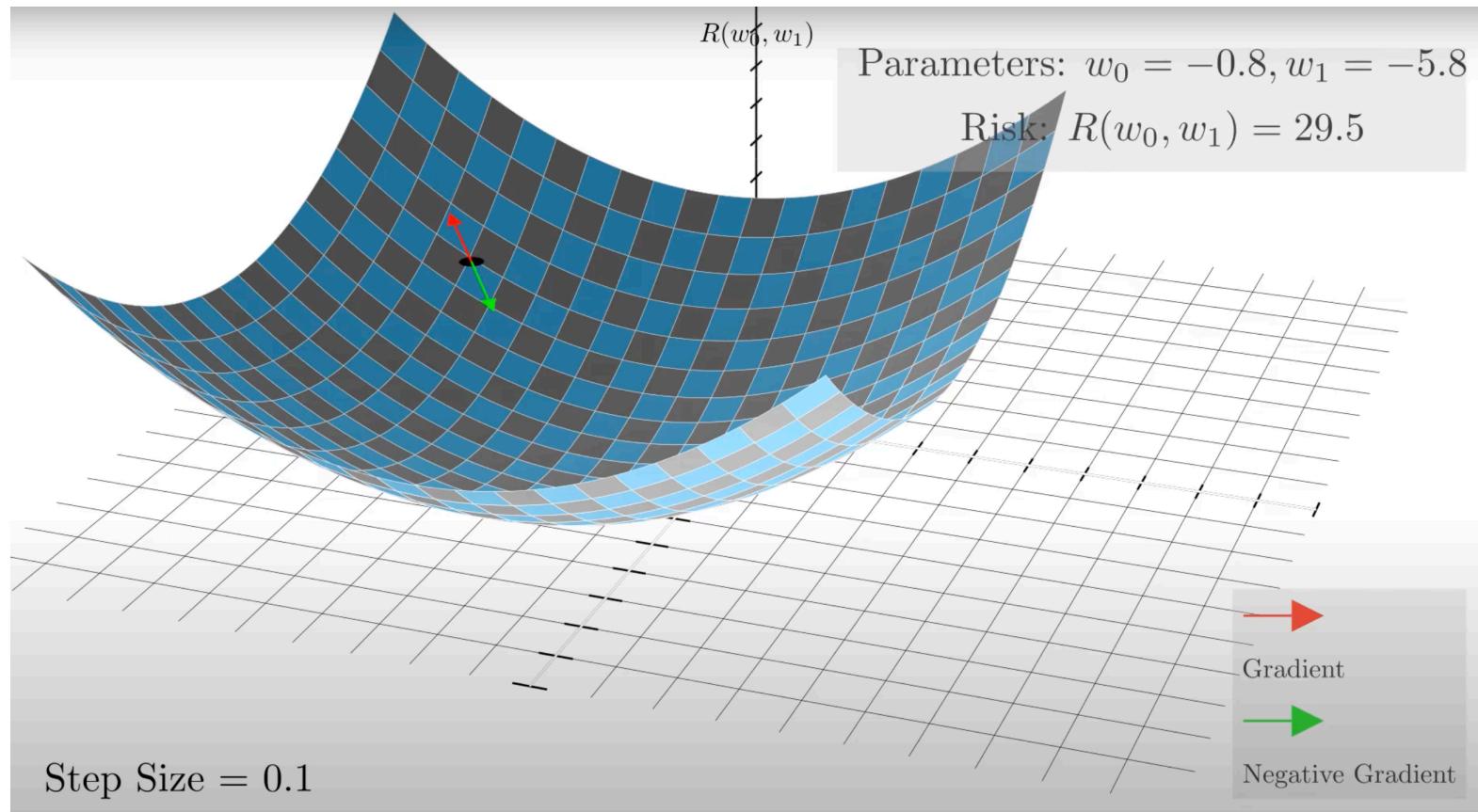
$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$

- 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?

- The Midterm Exam is tomorrow, in this room!
- In Homework 5, you'll see a few questions involving today's material:
 - A question about convexity.
 - A question about implementing gradient descent to find optimal parameters for a model that is **not linear in its parameters**.
- On Monday, we'll start talking about probability.
 - Homework 5 will have a probability problem taken from a past DSC 10 exam, to help you refresh.