

# **Gradient Descent**

# Learning Objectives

gradient descent = optimization algv.

- Describe what optimization is, ....
- ... and describe how gradient descent solves optimization problems
- Discuss what can go wrong when doing gradient descent
- Carry out gradient descent manually in Python

Learning Rate  
 $\eta$  (eta)

# \_ Optimization

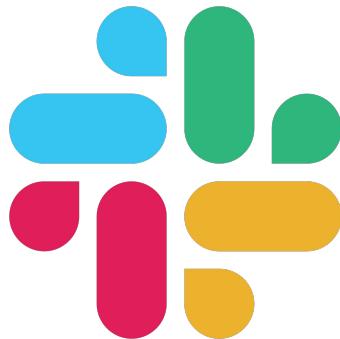
**Define: Optimization**  $\rightarrow$  ↳ use function / env  
                                  ↳ minimize

The term **optimization** refers to computing either the minimum or maximum of an **objective function**, maybe subject to some **constraints**. ]

↙ cust fn: regression  $\rightarrow$  MSE / MAE  
↙ loss fn: classification  $\rightarrow$  entropy .

There are many ways of solving an optimization problem, and you may already be familiar with one of them: by computing the **derivative** of a function and setting it to zero.

# Define: Derivative



What is a derivative?

# Define: Derivative

A **derivative** of a function at a point is the value of the slope of the tangent line at that point.

$$\frac{dy}{dx}$$



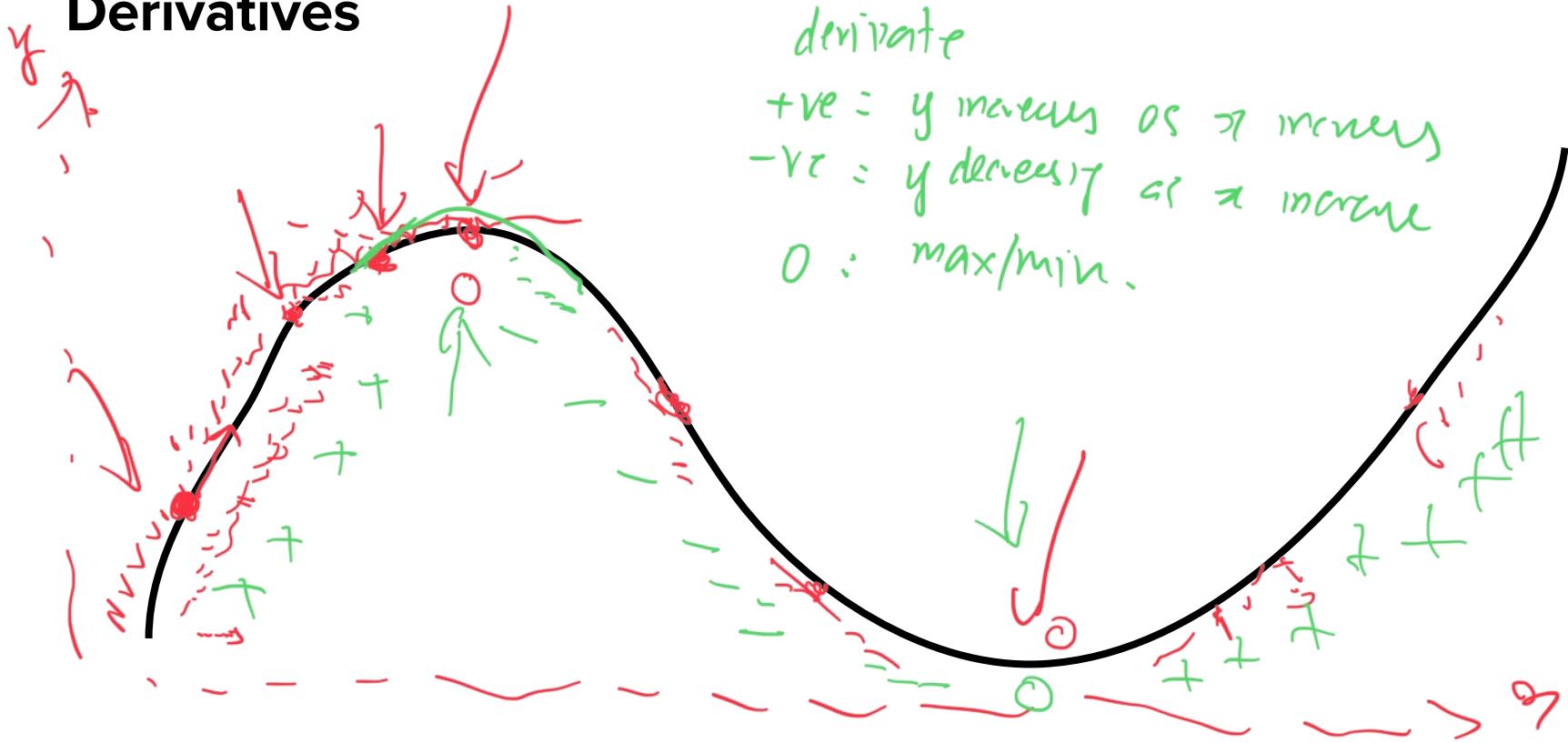
Training  
Loss



$$y = \underline{B_0}x + \underline{B_1}x^2$$

optimization  
algorithm  
that helps you  
find the optimal  
parameters for your model  
given loss function. GA

# Derivatives



# Derivatives

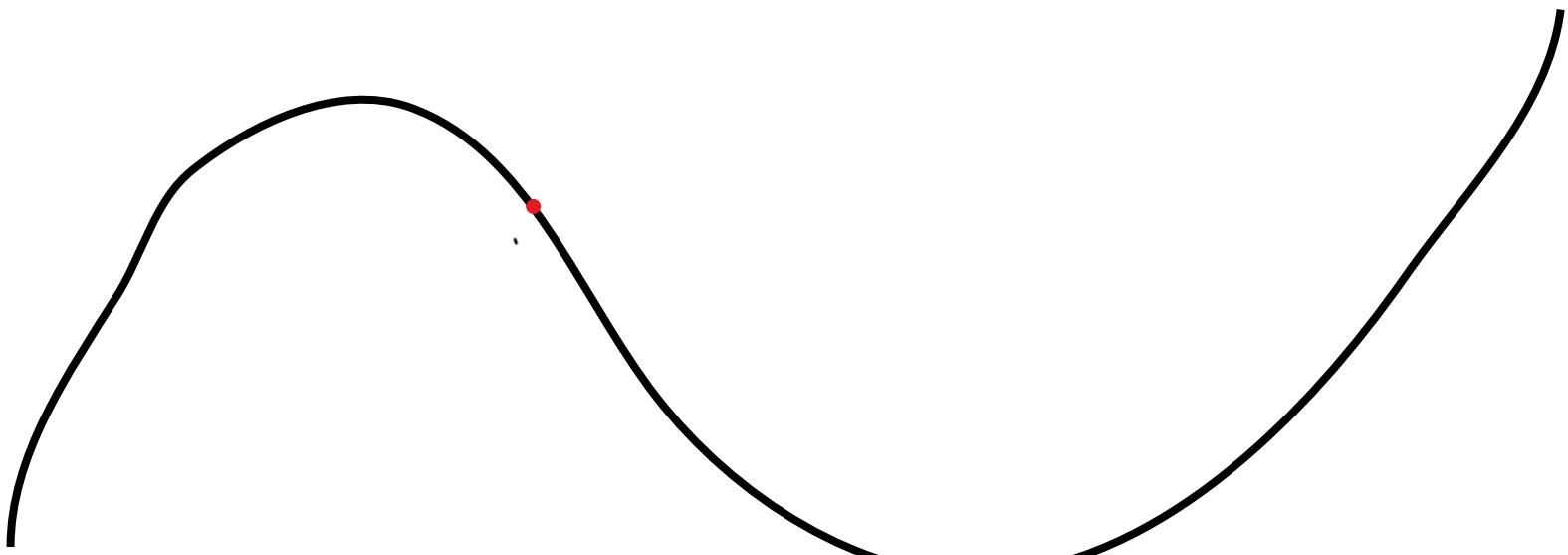
derivative = slope.

↳ represents steepness of the curve  
at a given point

Sign + absolute value



# Derivatives



# Derivatives

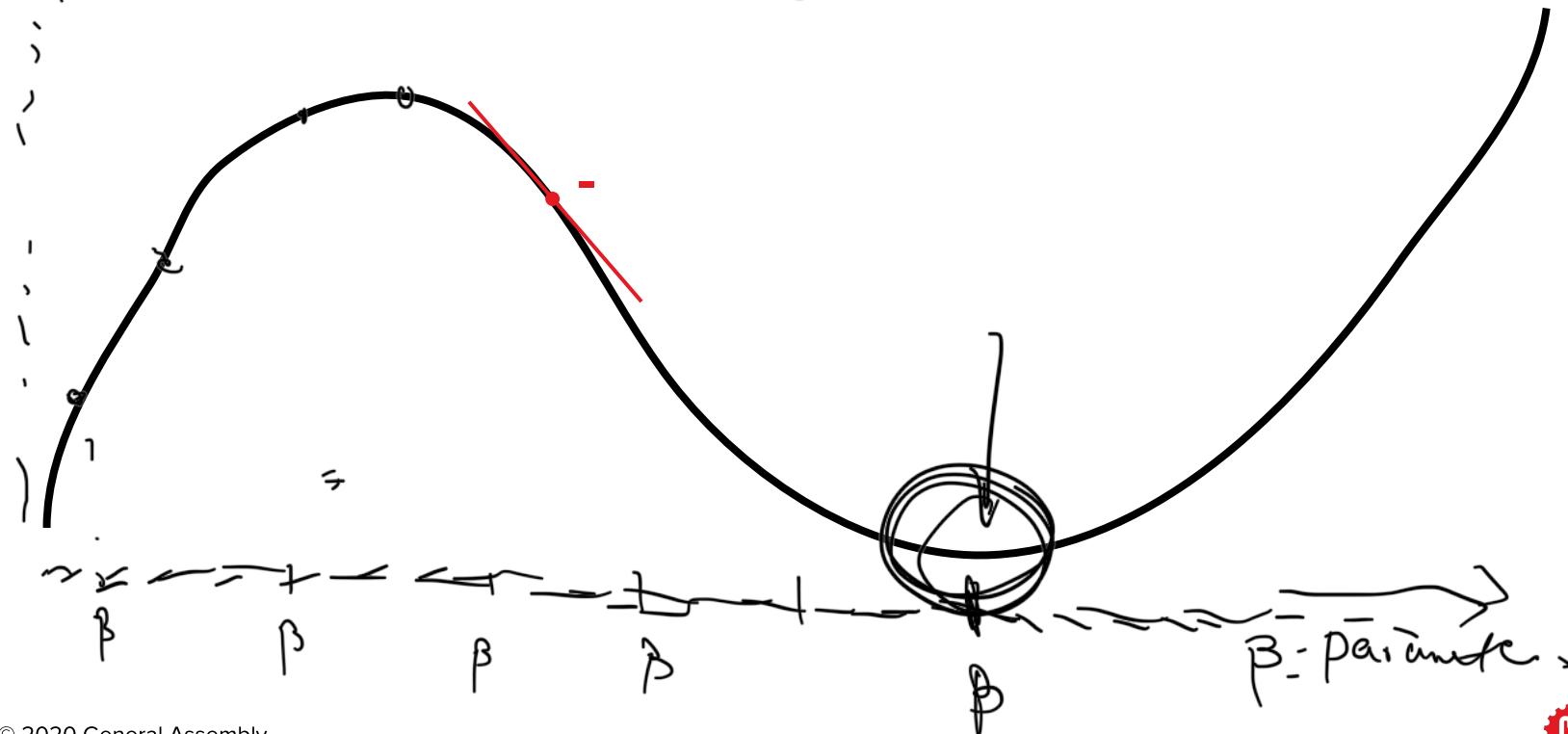
Loss fn



$$y = \beta x$$

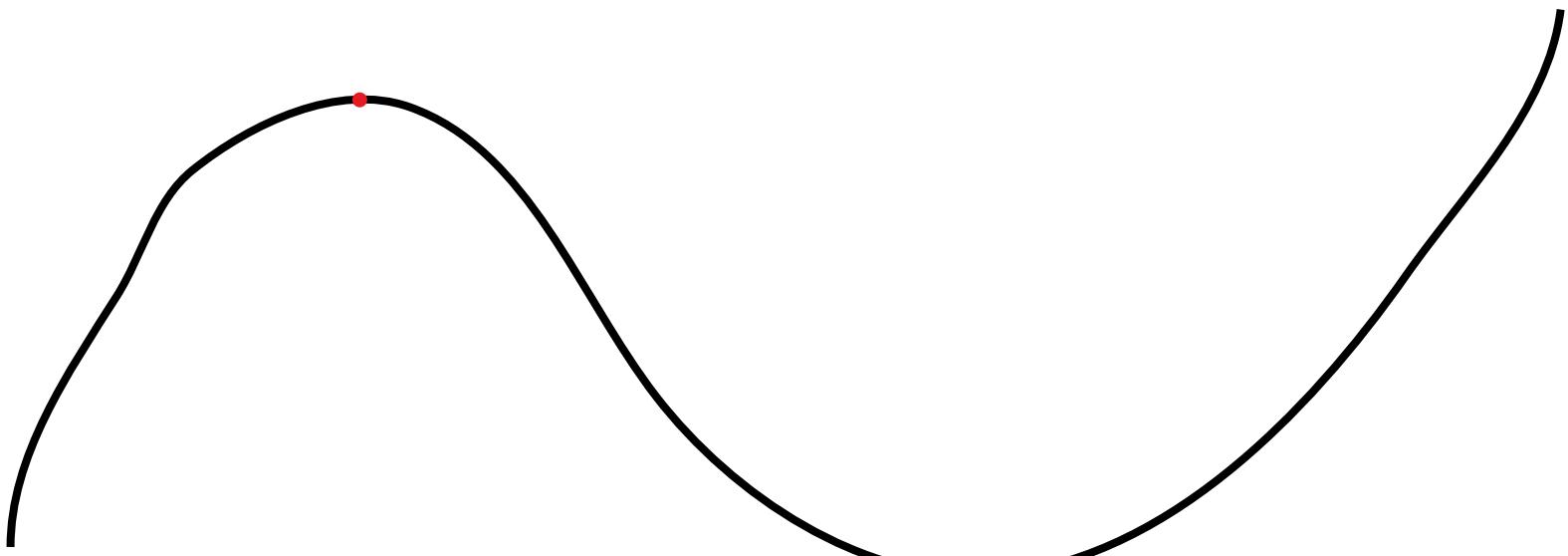


$$y = \beta_1 x + \beta_2 x^2$$

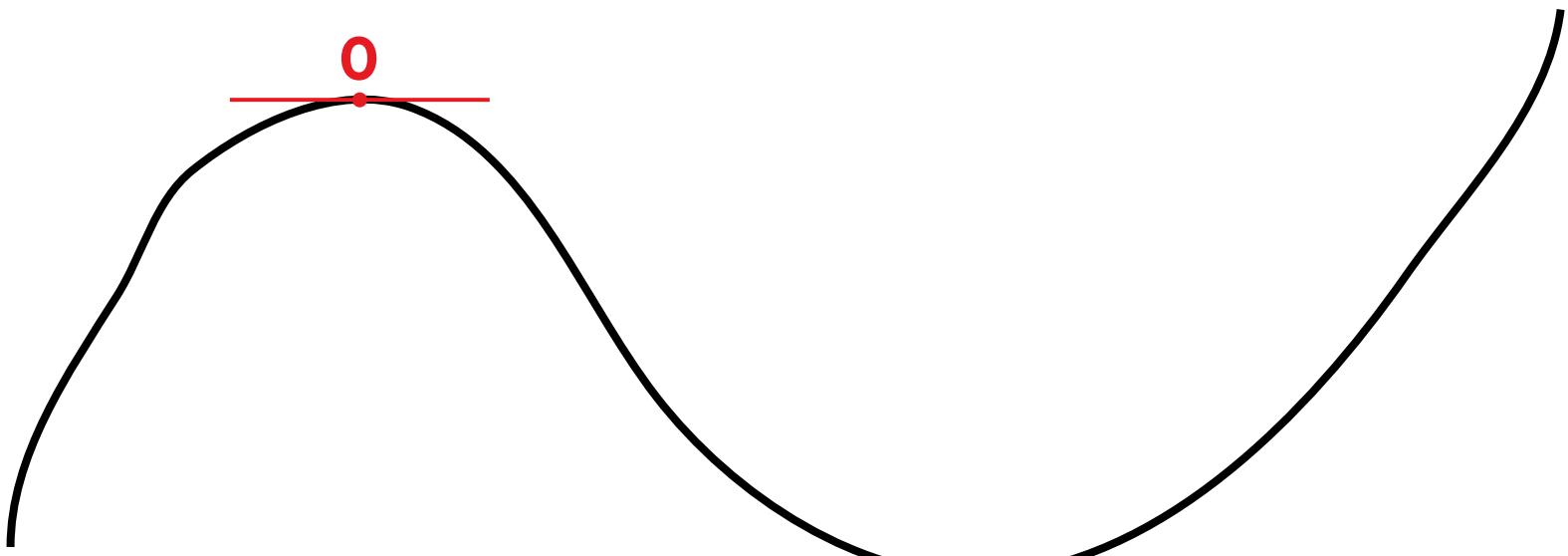


---

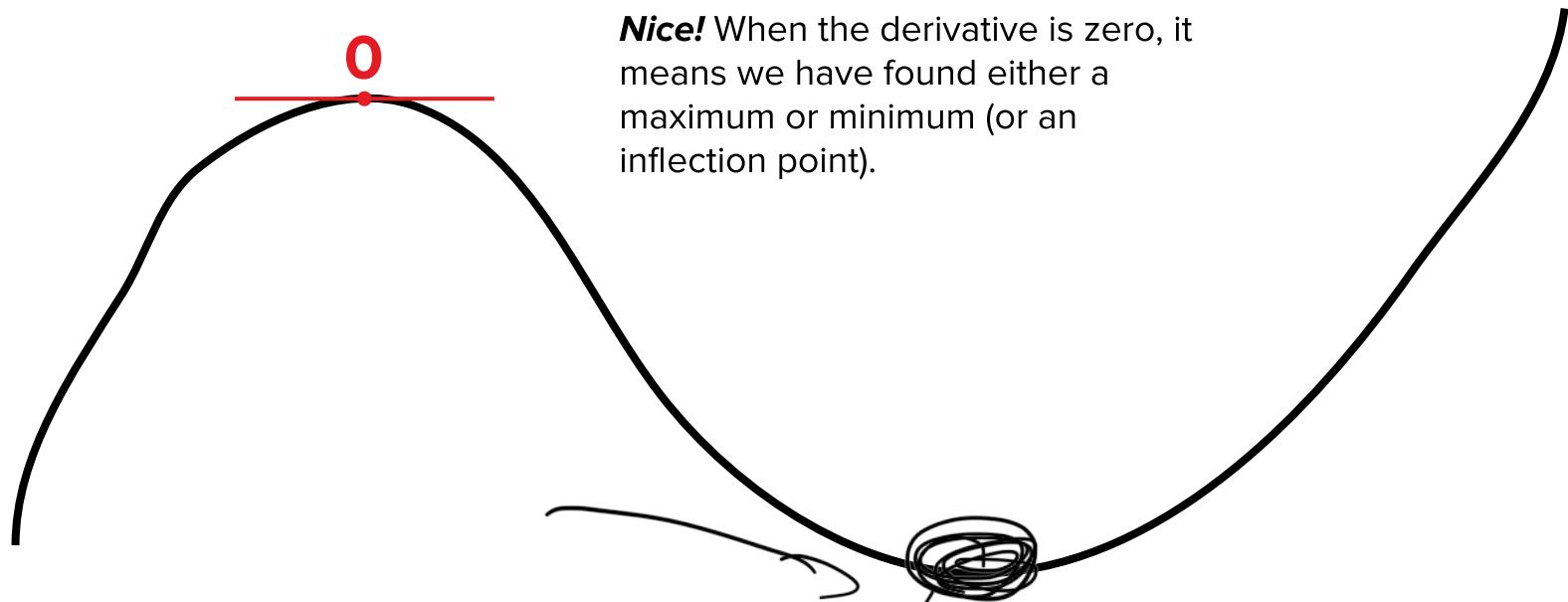
# Derivatives



# Derivatives



# Derivatives



# Optimizing the “old fashioned way”

In calculus class, you might have optimized a function by computing its derivative, setting it to zero, and then solving for x.

~~y~~

y

$$\frac{dy}{dx} = 0 \Rightarrow x ?$$

# Optimizing the “old fashioned way”

In calculus class, you might have optimized a function by computing its derivative, setting it to zero, and then solving for  $x$ .

***Today's lesson is about what happens when you can't do that!***



# Optimizing for data science

(Almost) all machine learning algorithms involve some sort of optimization problem. With one notable exception (OLS), they are all unsolvable the “old fashioned way”.

Instead, we'll need to use **gradient descent** to optimize **loss functions**.

- 
- gradient descent : recipe to train a model
- ✓ (1) loss function / cost fn / objective fn
  - ✓ (2) training data
  - ✓ (3) model architecture. ↪ SVM neural network  
↪ what level.
- The notes are handwritten in black ink. A large bracket on the right side groups the three items: "loss function", "training data", and "model architecture". To the left of this bracket, there are three checkmarks and three arrows pointing towards the bracket. Above the bracket, the text "gradient descent : recipe to train a model" is written in a larger, more formal font. Below the bracket, the three items are listed with checkmarks and arrows. The word "SVM" is written above "neural network".

## Define: Loss function

$$\beta_{t+1} = \beta_t - \eta \nabla_{\beta_t} \text{loss}_{f_n}$$



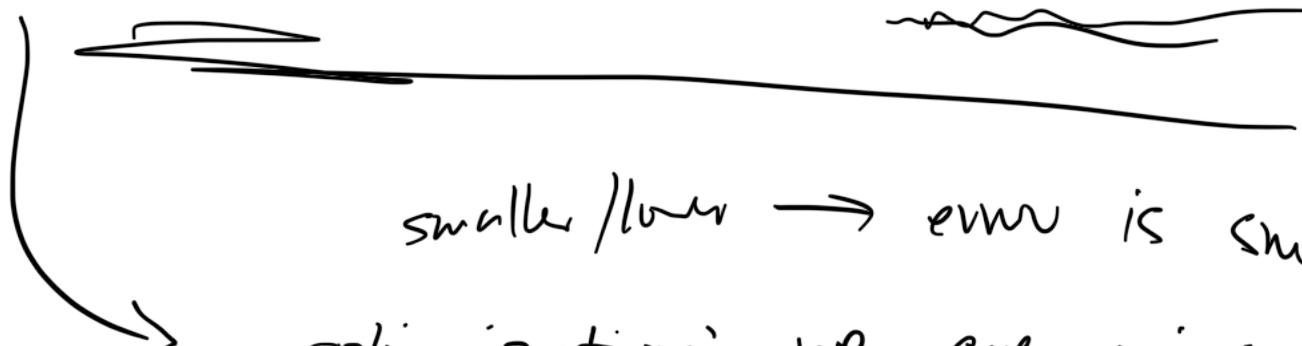
In your own words, what is a loss function?



Name all the loss functions you know.

## Define: Loss function

A **loss function** is a function that describes the difference between  $y$  and  $\hat{y}$ . In order for something to be a loss function, **lower values must be better**.



smaller / lower  $\rightarrow$  err is small.

optimization: we are minimizing  
the loss fn

## Define: Loss function

A **loss function** is a function that describes the difference between  $y$  and  $\hat{y}$ . In order for something to be a loss function, **lower values must be better**.

This means that  $R^2$  (for regression) and accuracy (for classification) don't count as loss functions!



# OLS: Minimizing the MSE loss function

The goal of OLS is to find the value of  $b$  that minimizes the MSE:

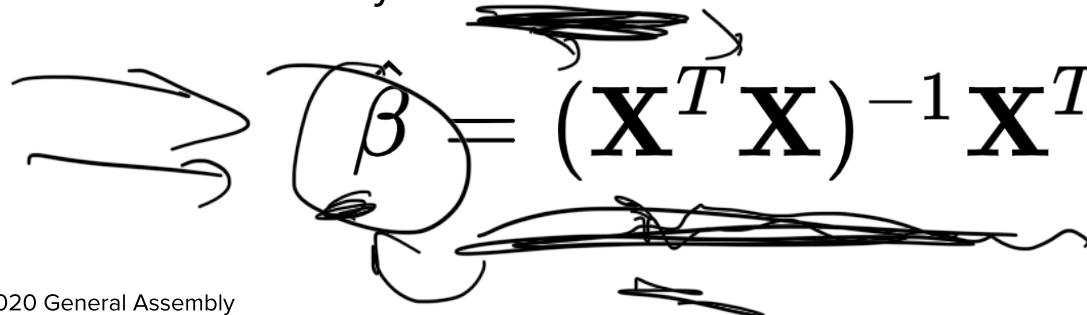
$$\text{minimize} \quad \sum (y_i - \mathbf{x}_i^T \beta)^2$$
$$(y_i - \beta^T x_i)^2$$

# Minimizing the MSE loss function

The goal of OLS is to find the value of  $b$  that minimizes the MSE:

$$\text{minimize} \quad \sum (y_i - \mathbf{x}_i^T \beta)^2$$

This is one of the only problems interesting to data scientists that we can solve the “old fashioned way” in closed form:


$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

# Minimizing the LASSO loss function

The LASSO, on the other hand, is harder. We can't solve this with simple calculus:

~~minimize~~  $\sum (y_i - \mathbf{x}_i^T \beta)^2 + \lambda \sum |\beta_j|$  ↴ we want  $\beta$  to be small.

The diagram illustrates the LASSO loss function, which is the sum of the squared errors and a multiple of the sum of the absolute values of the coefficients. The wavy line represents the squared error term, and the jagged line represents the absolute value regularization term. The goal is to find a balance where the coefficients are small enough to fit the data well but not zeroed out entirely.

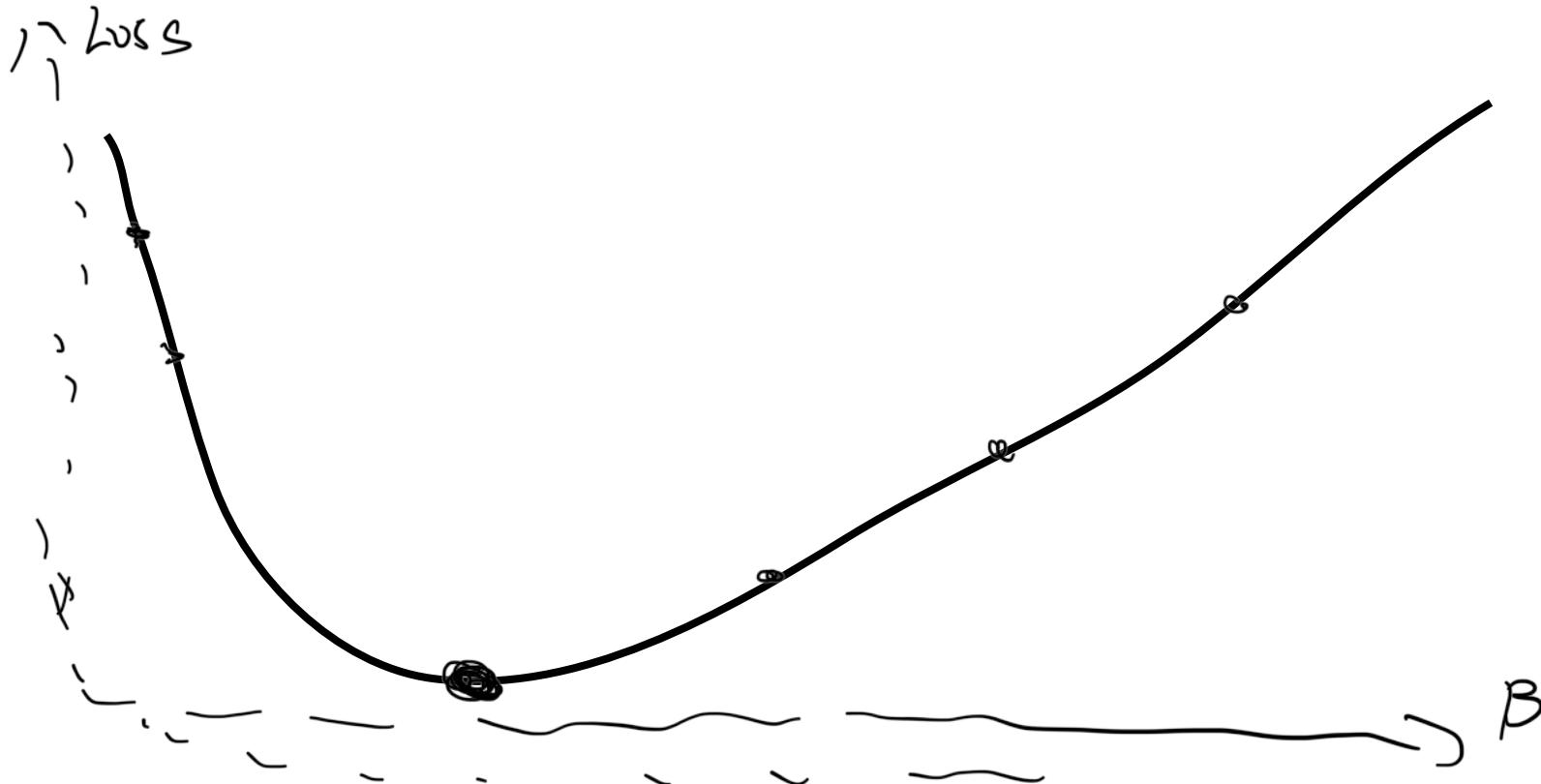
## Minimizing the log-loss function

Still more difficult is the **log-loss function** (also called **binary crossentropy**) for finding the optimal betas for logistic regression:

$$\text{minimize} \quad \sum [-y_i \log \hat{y}_i - (1 - y_i) \log(1 - \hat{y}_i)]$$

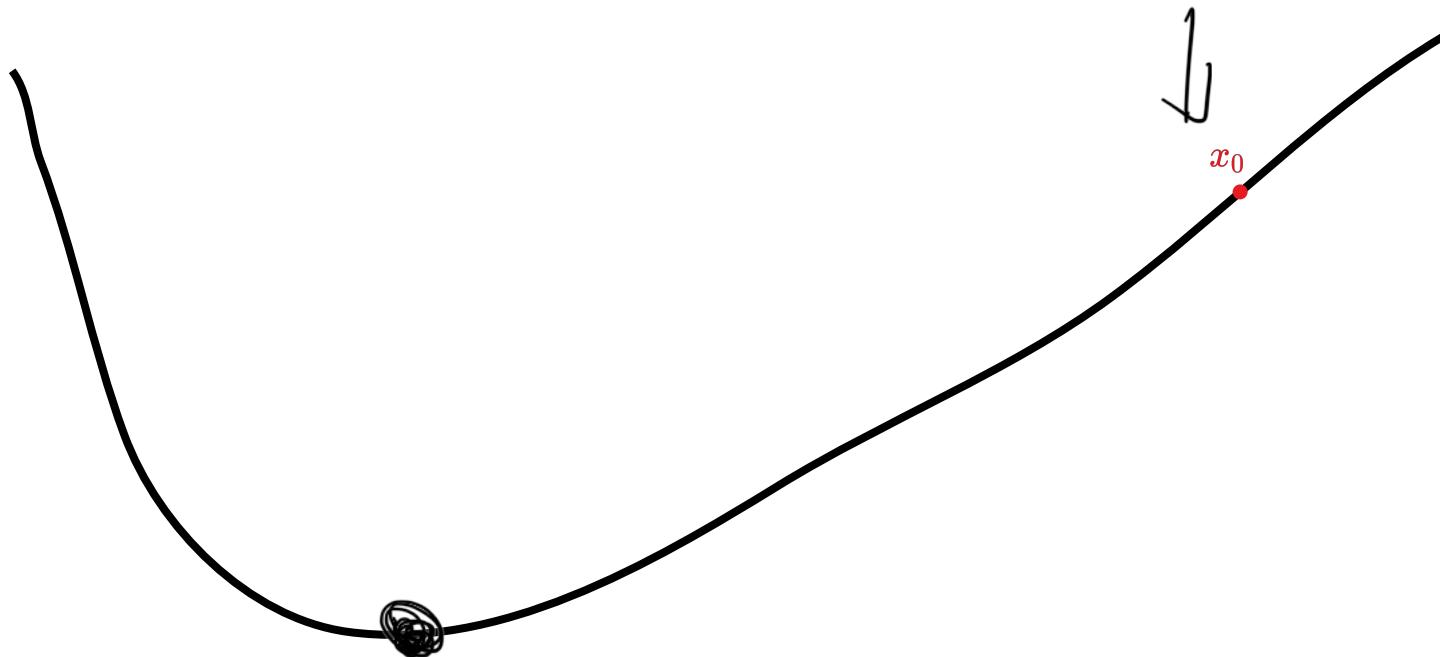
# - Optimizing loss functions with gradient descent

# Gradient Descent



---

## Step 0: Start at a random point

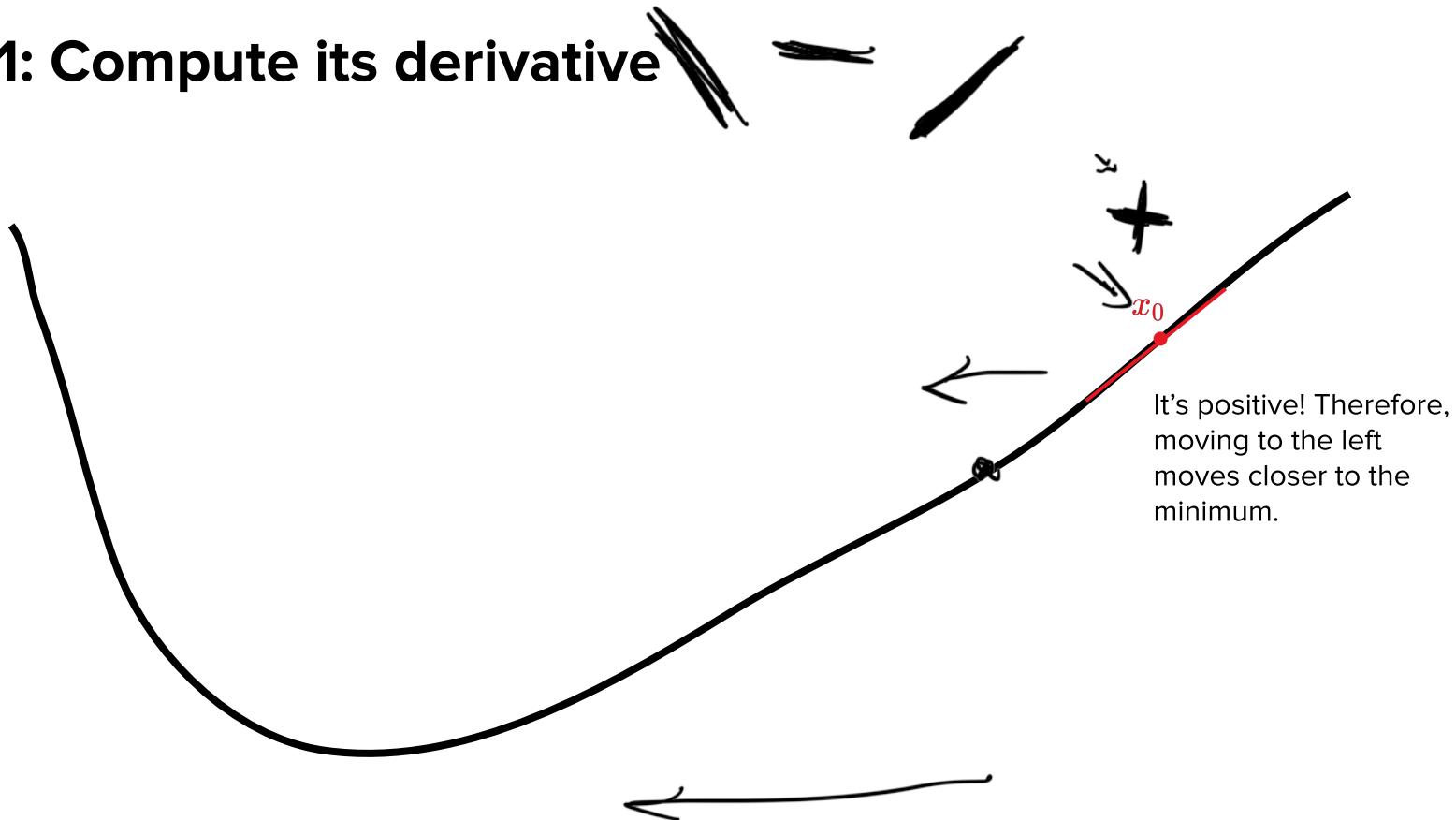


---

## Step 1: Compute its derivative

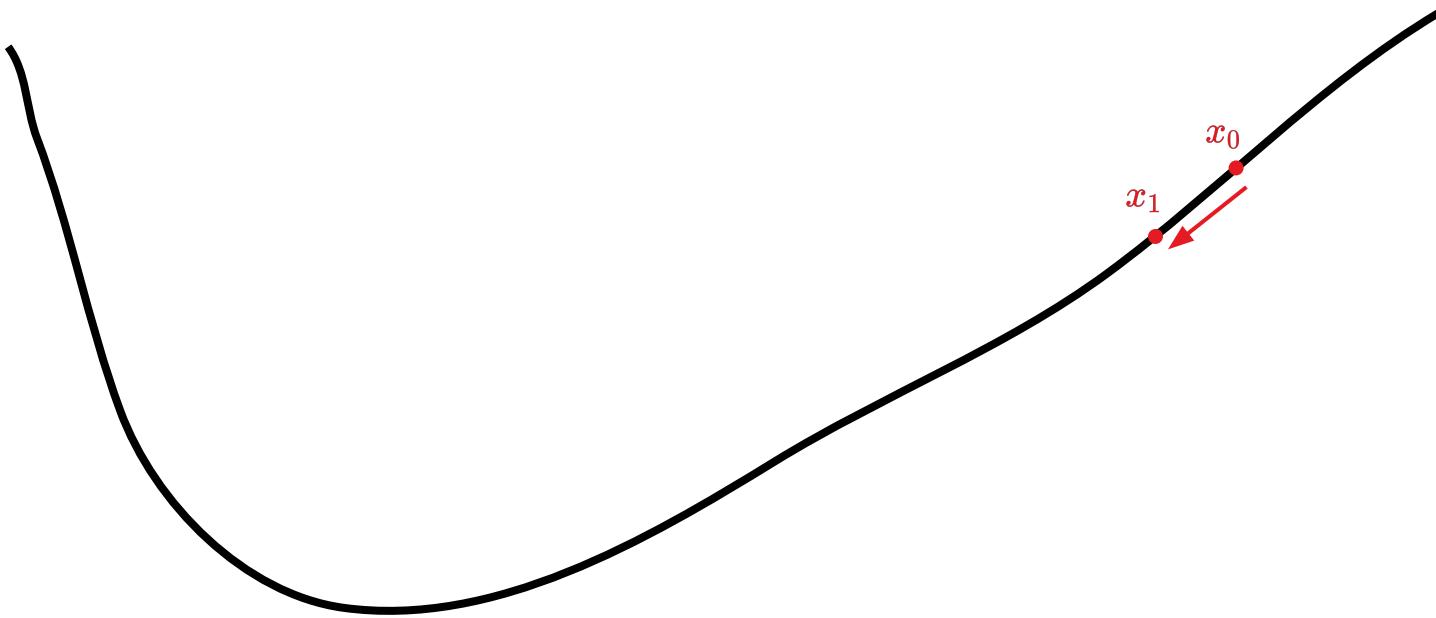


# Step 1: Compute its derivative



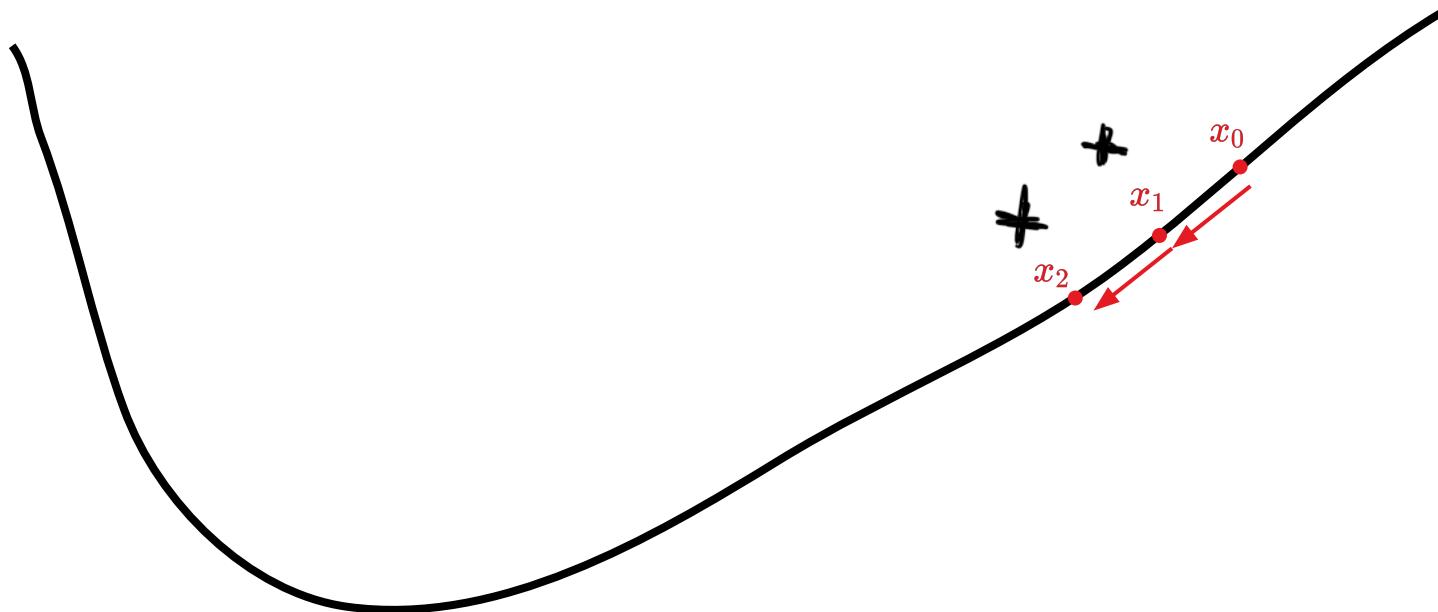
---

## Step 2: Move in the appropriate direction



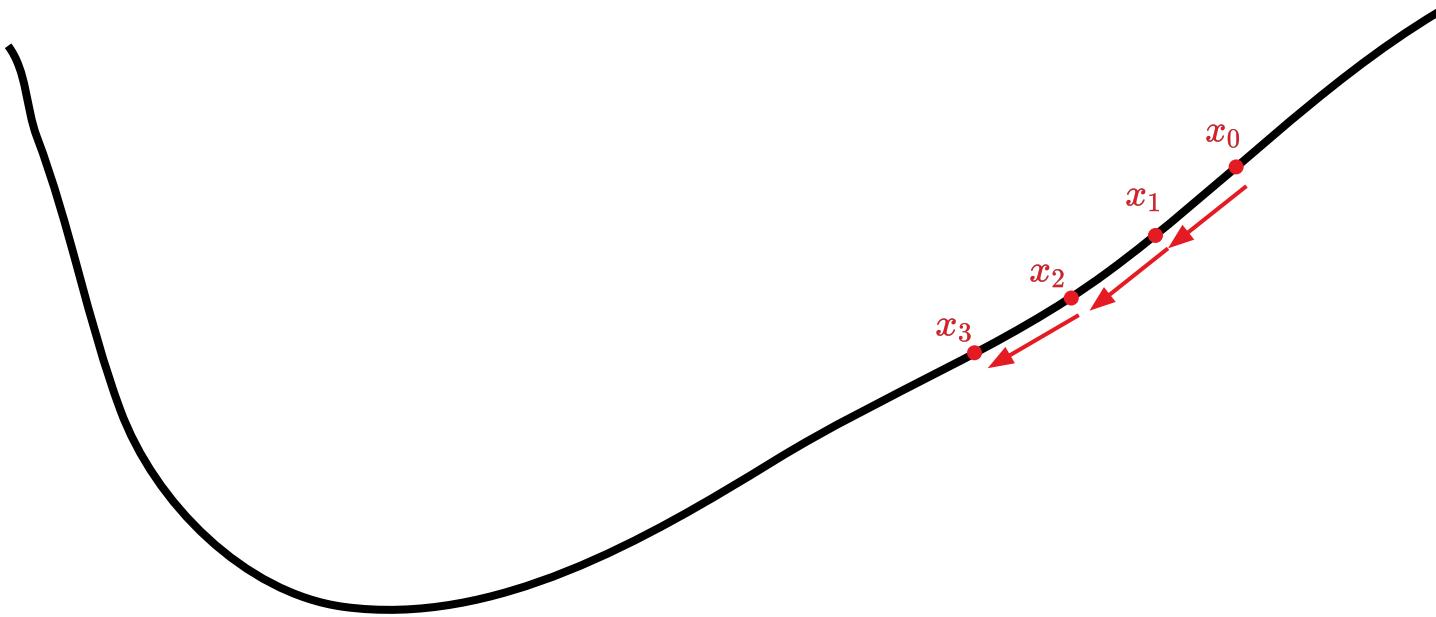
---

## Step 3: Keep repeating steps 1 and 2

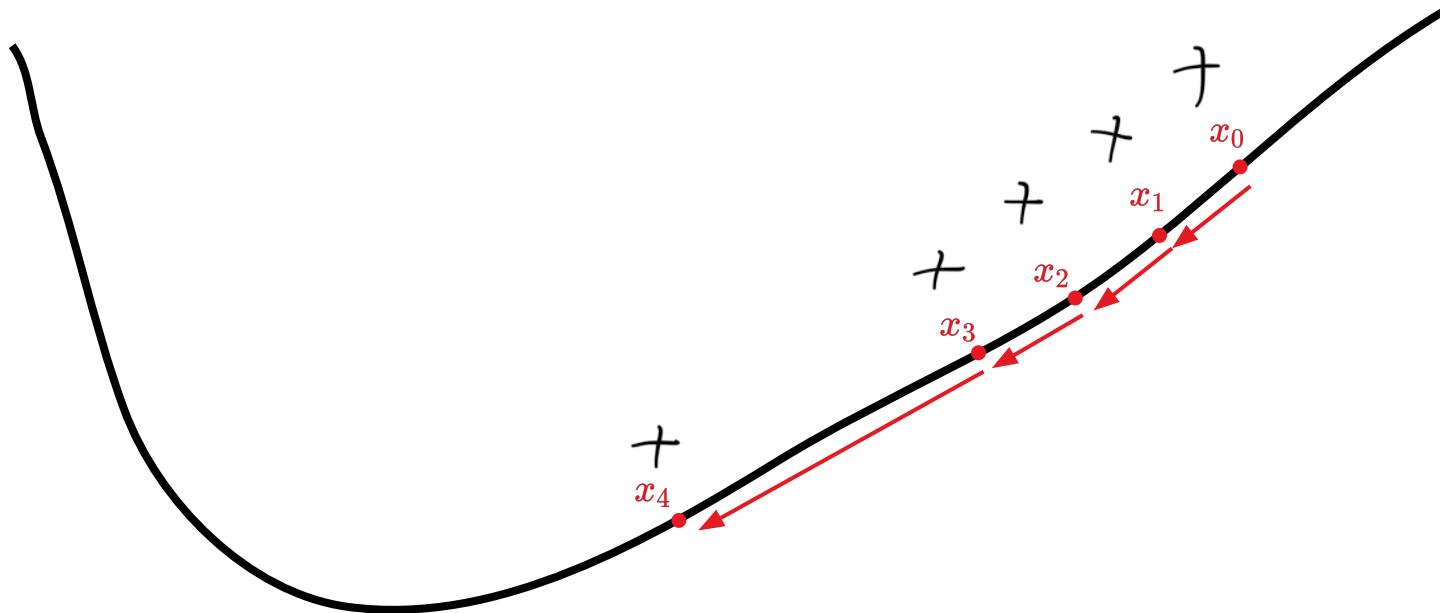


---

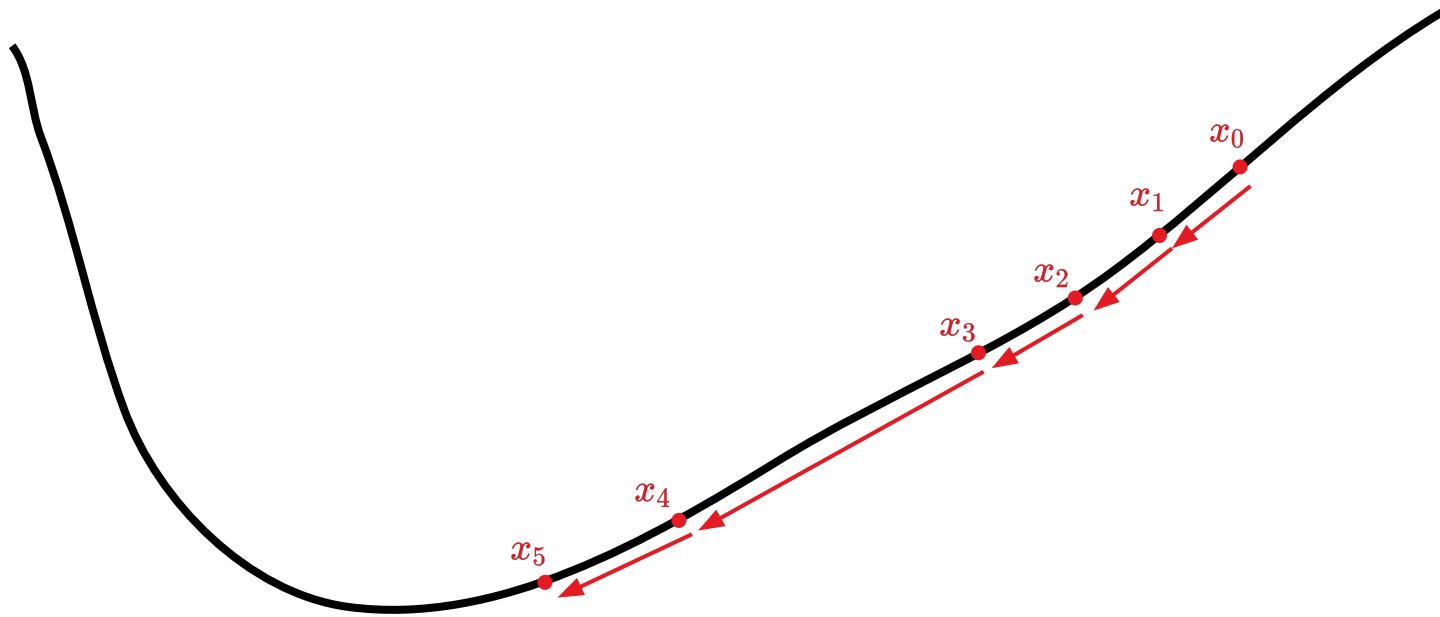
## Step 3: Keep repeating steps 1 and 2



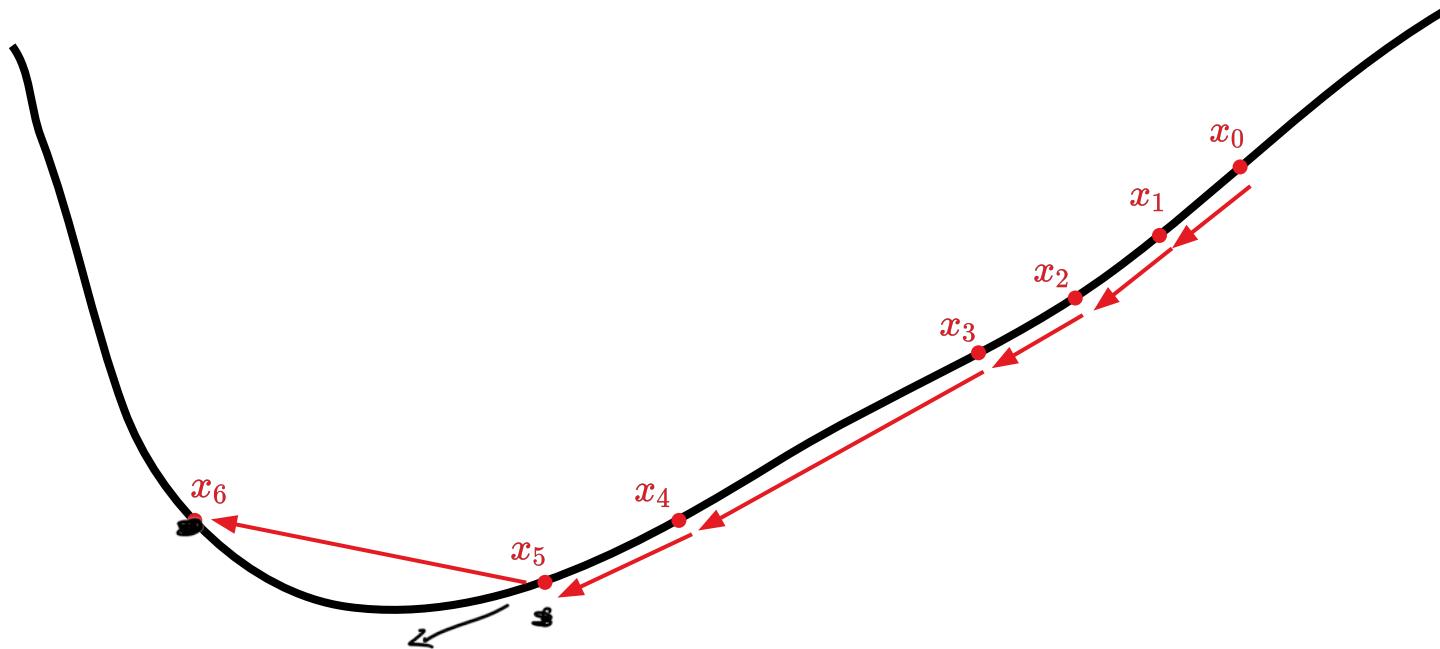
## Step 3: Keep repeating steps 1 and 2



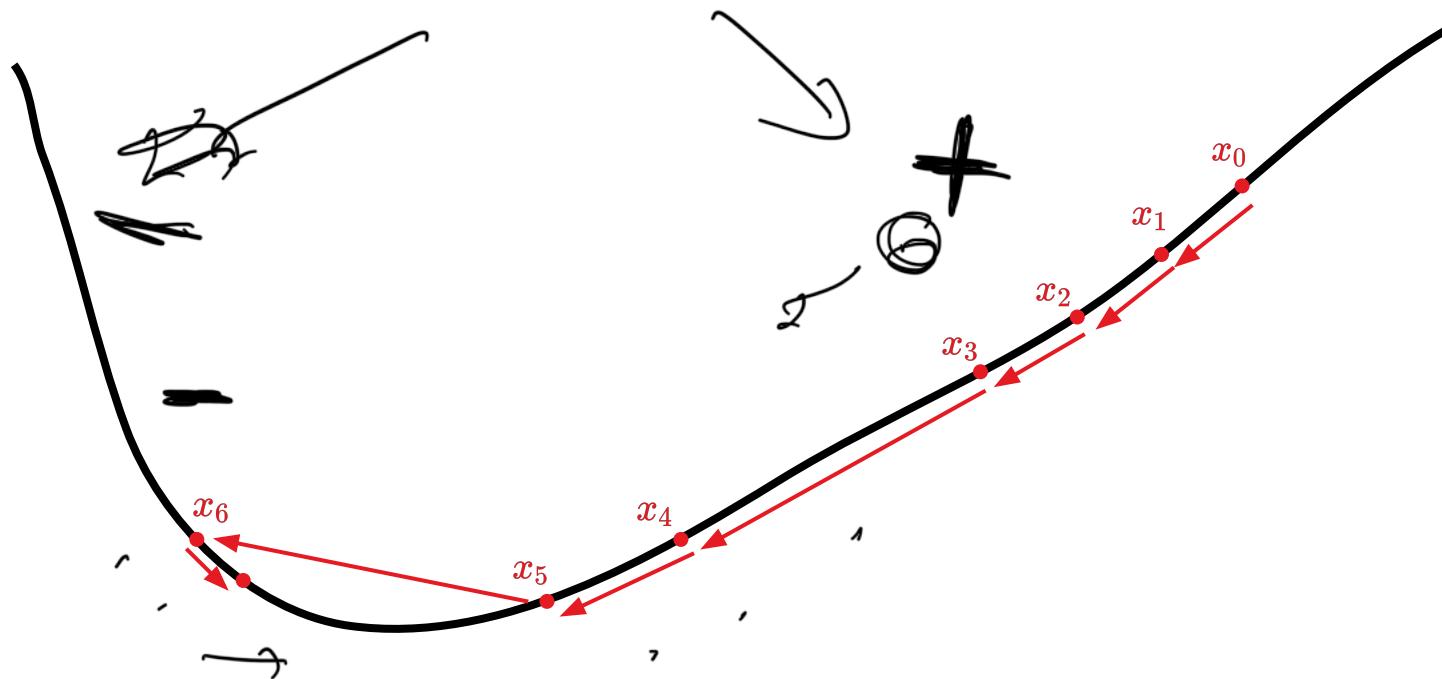
## Step 3: Keep repeating steps 1 and 2



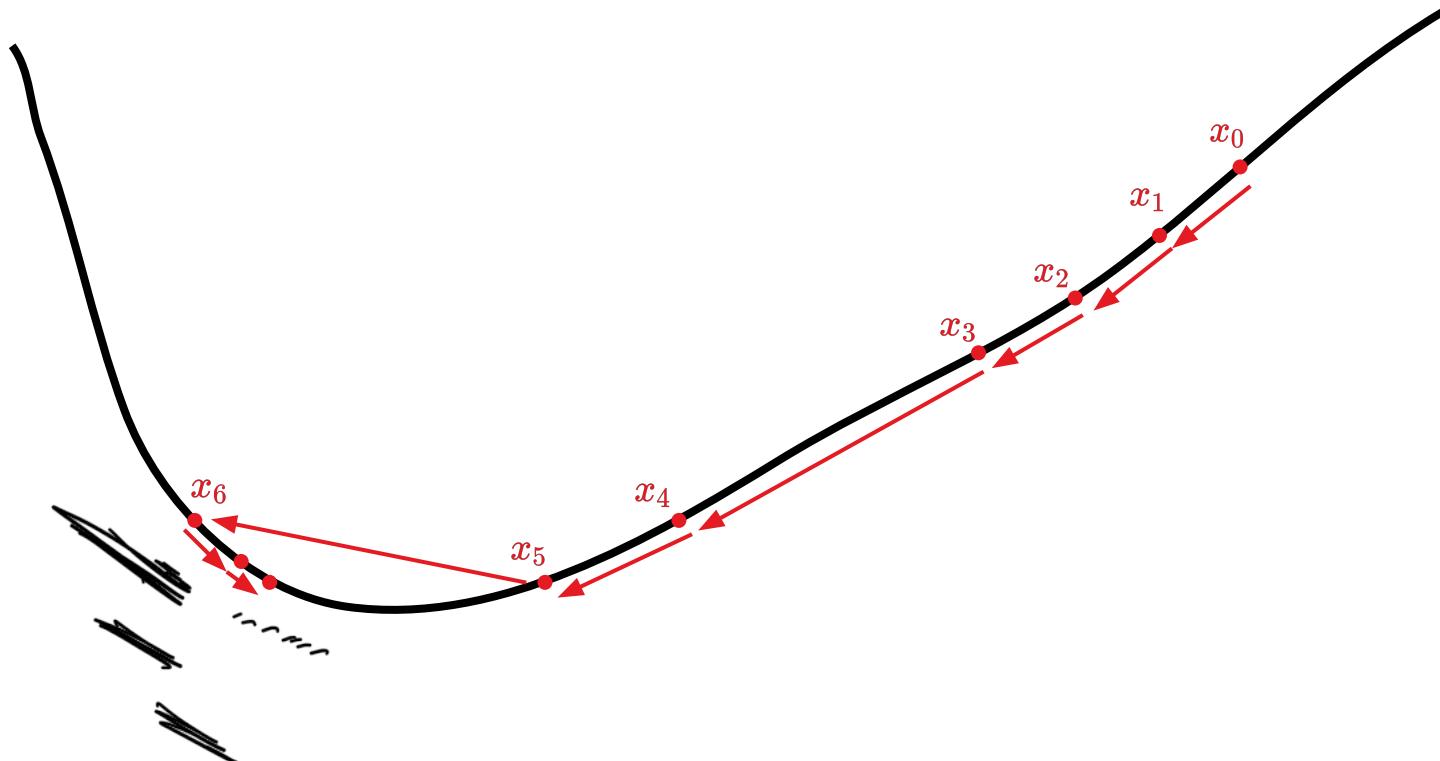
## Step 3: Keep repeating steps 1 and 2



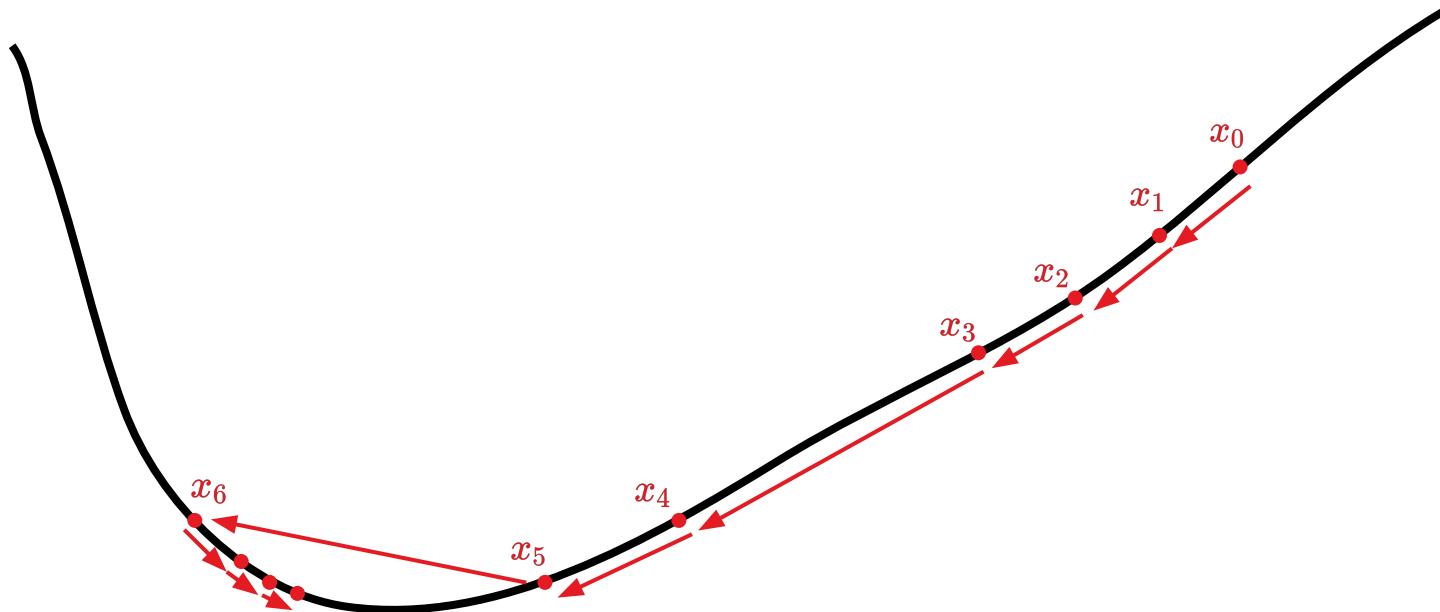
## Step 3: Keep repeating steps 1 and 2



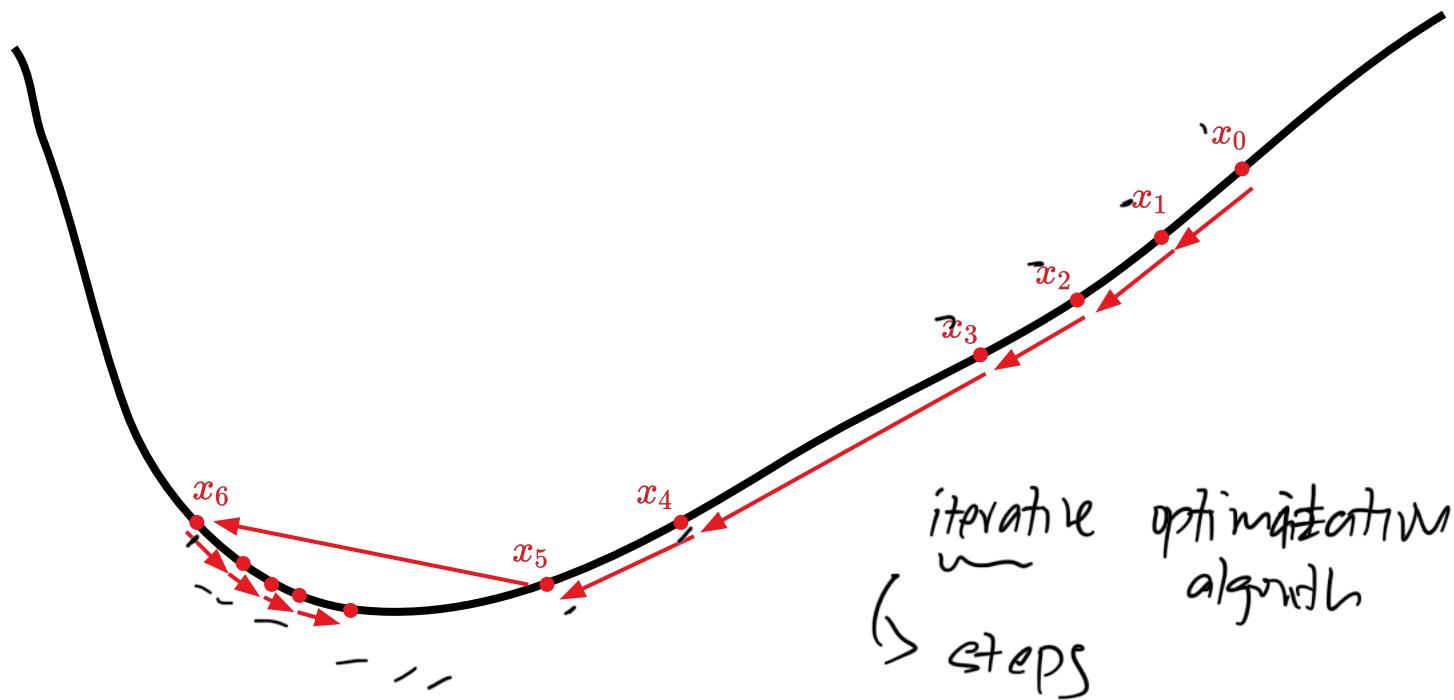
## Step 3: Keep repeating steps 1 and 2



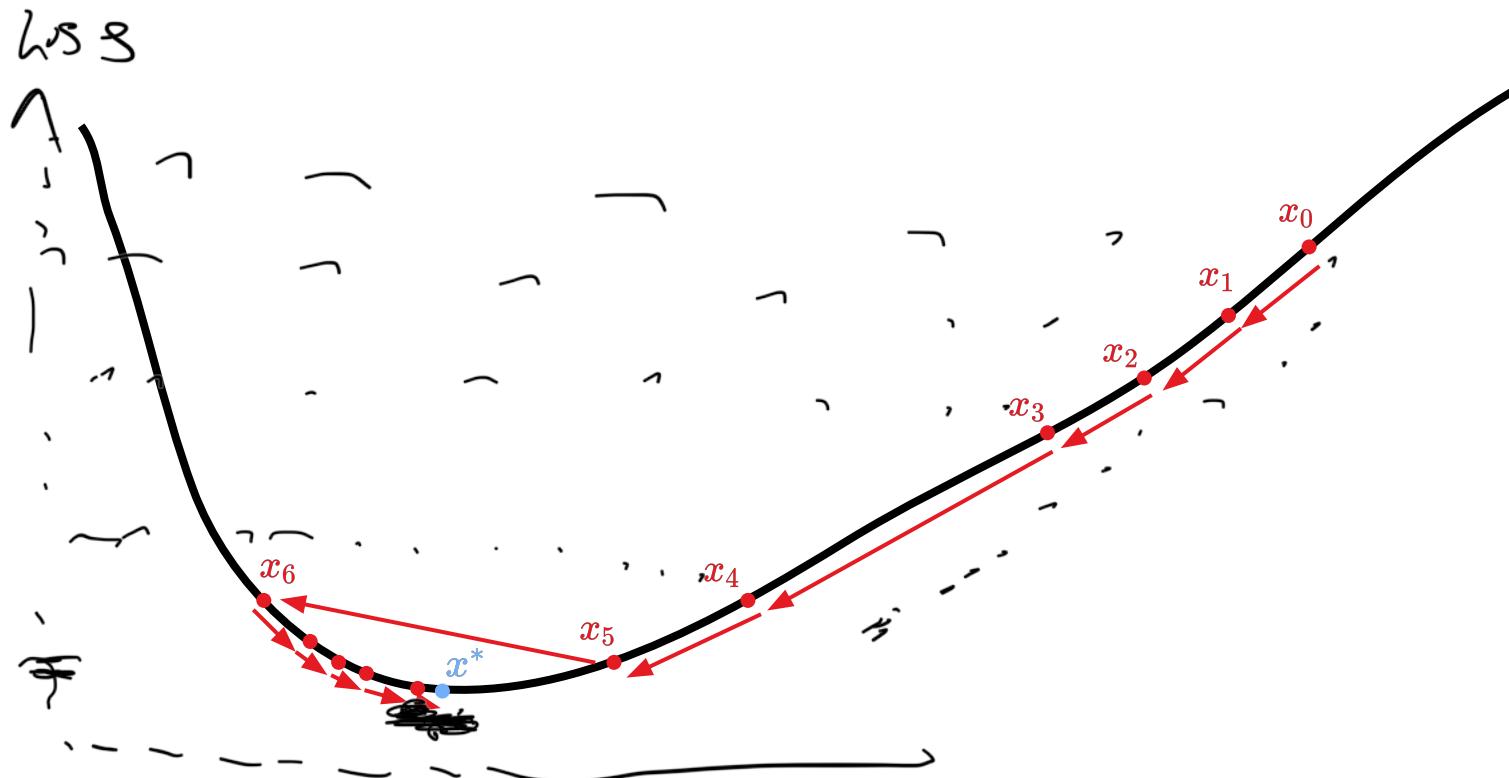
## Step 3: Keep repeating steps 1 and 2



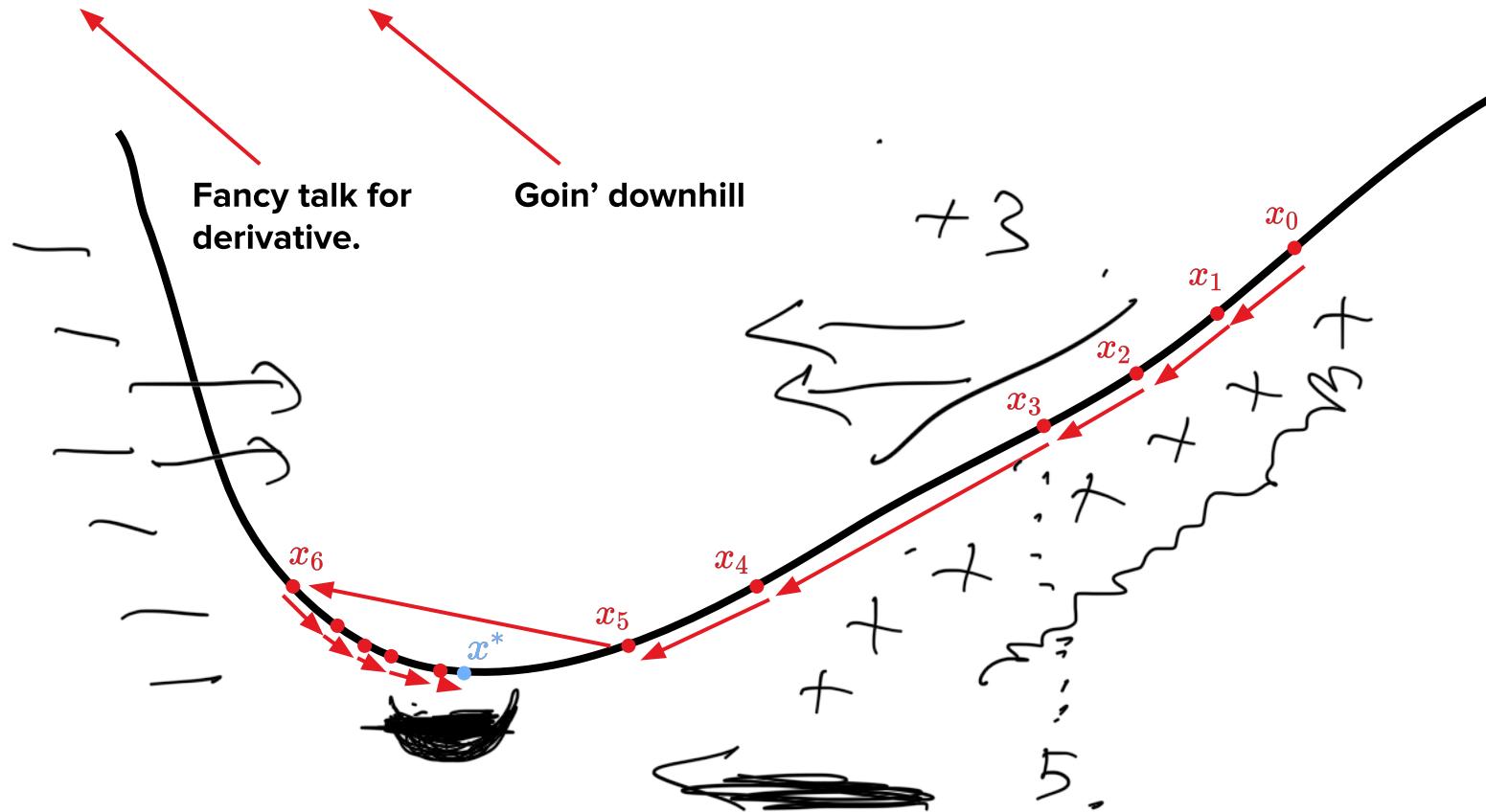
## Step 3: Keep repeating steps 1 and 2



## Step 4: Stop when your move is small enough



# Gradient Descent



# Gradient Descent

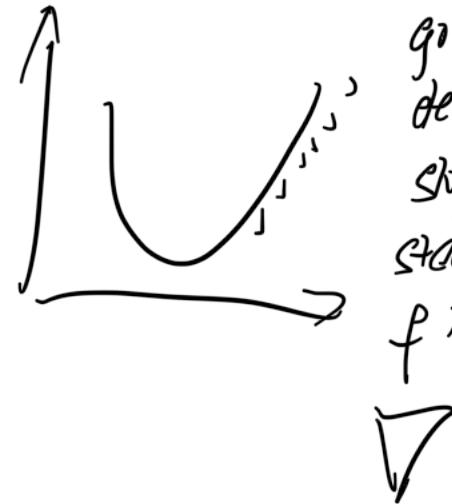
So how does it work, mathematically? There are a few considerations:

- We'll be “updating” our current position over and over again
  - We need to move in the opposite direction of the derivative
  - We need to decide how far to move at each step
- 

# Gradient Descent

$$\underline{x_{k+1}} \leftarrow \underline{x_k}$$

New position                      Old position



gradient /  
derivative /  
slope /  
steepness.  
 $f'$   
 $\nabla$

# Gradient Descent

$$x_{k+1} \leftarrow x_k - f'(x_k)$$

*f(x)*  
*f'(x) =*

New position      Old position      ... the gradient

Opposite direction of...

# Gradient Descent

$$x_{k+1} \leftarrow x_k - \alpha f'(x_k)$$

New position

Old position

Opposite direction of...

... the gradient

Step size



$\eta$  (lota)  
learning rate



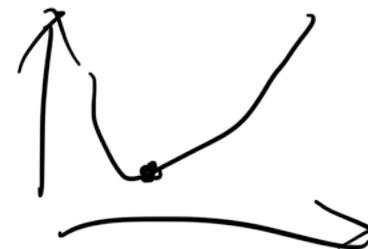
---

# What can go wrong?

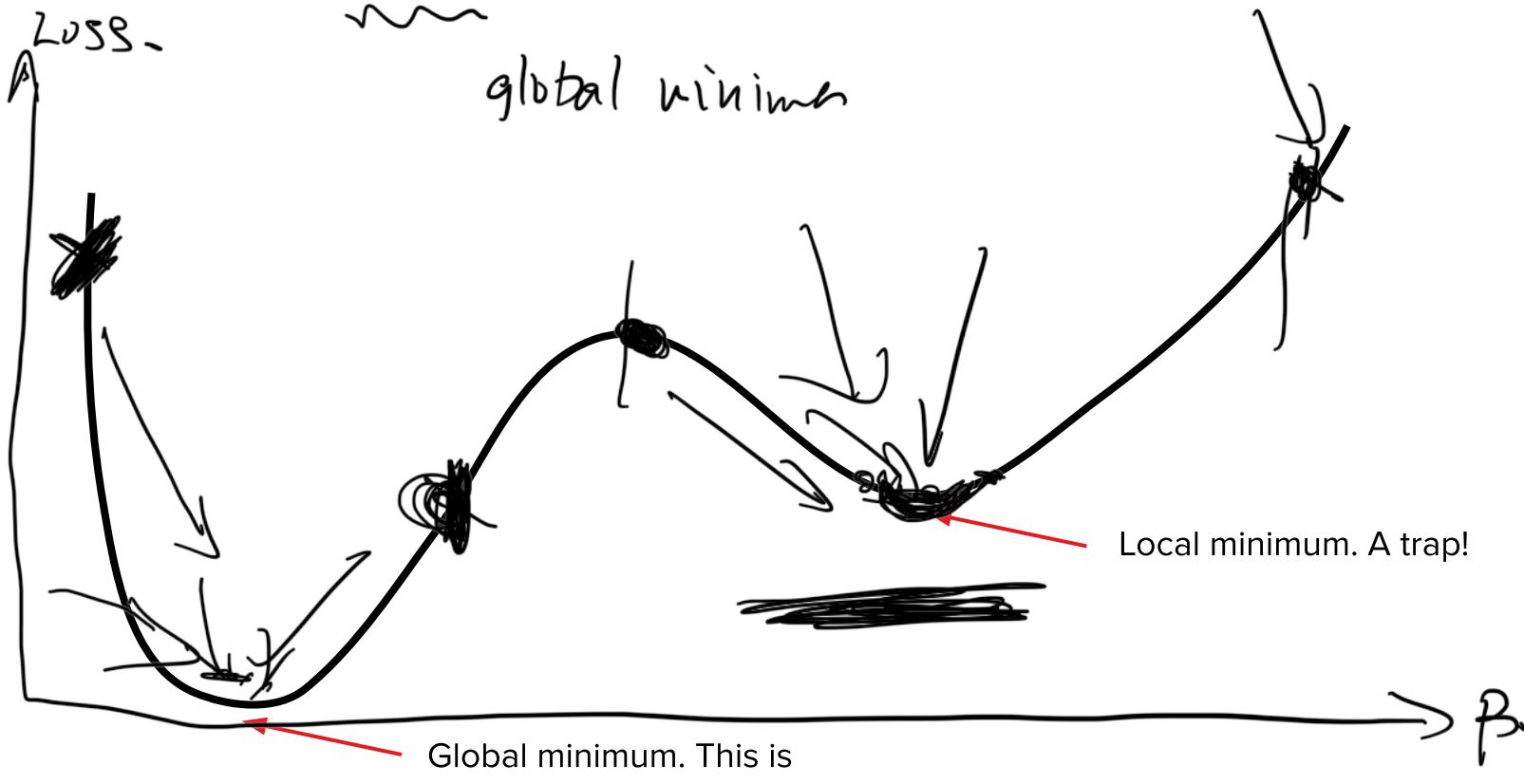
# What can go wrong?

While it's difficult to grasp at first, gradient descent is actually a very simple algorithm. There's really only one component that's difficult to calculate, and the rest is easy.

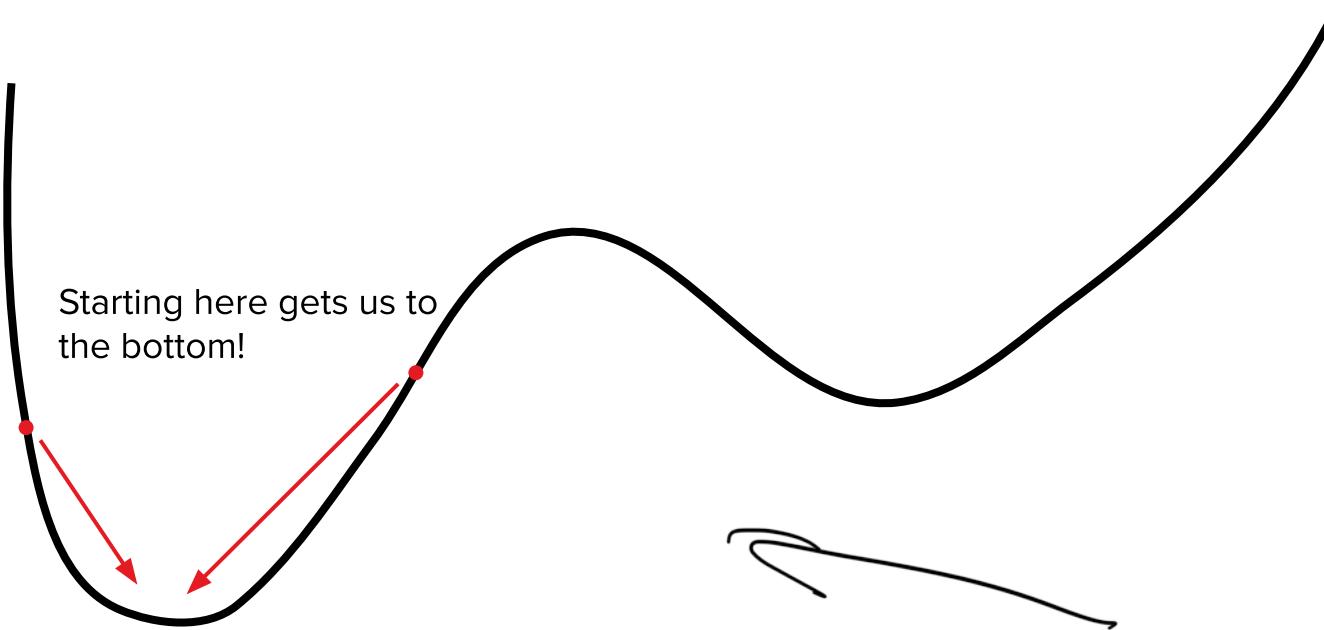
But if it were *really* easy, why make such a big fuss?



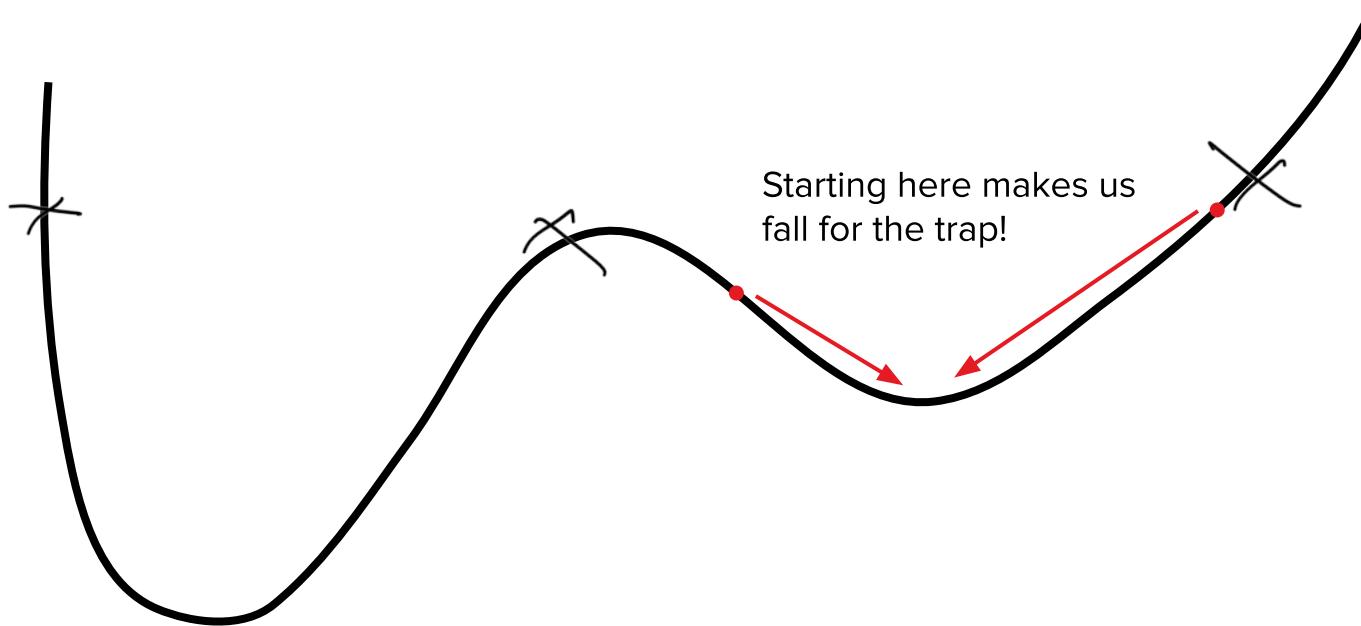
## Problem #1: Local Minima



## Problem #1: Local Minima



# Problem #1: Local Minima



## Problem #1: Local Minima

The solution to this problem is to verify that our objective function is convex.

A function is **convex** if it is perfectly U-shaped.

Want a more formal definition?



# Convexity

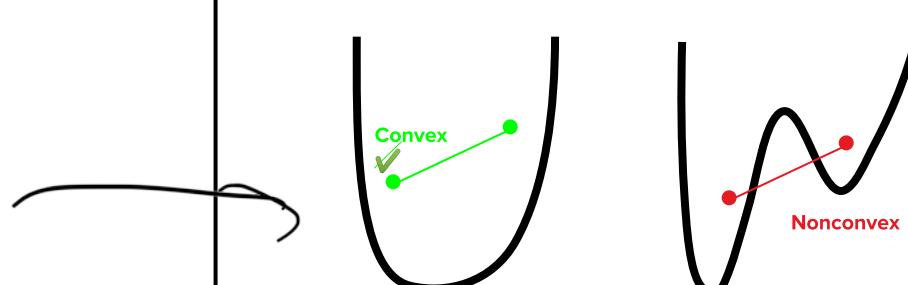
A function is **convex** if the two equivalent definitions are true:

**Definition #1**

$$f''(x) \geq 0 \text{ for all } x$$

**Definition #2**

The straight line between any two points above a function are entirely inside that function.

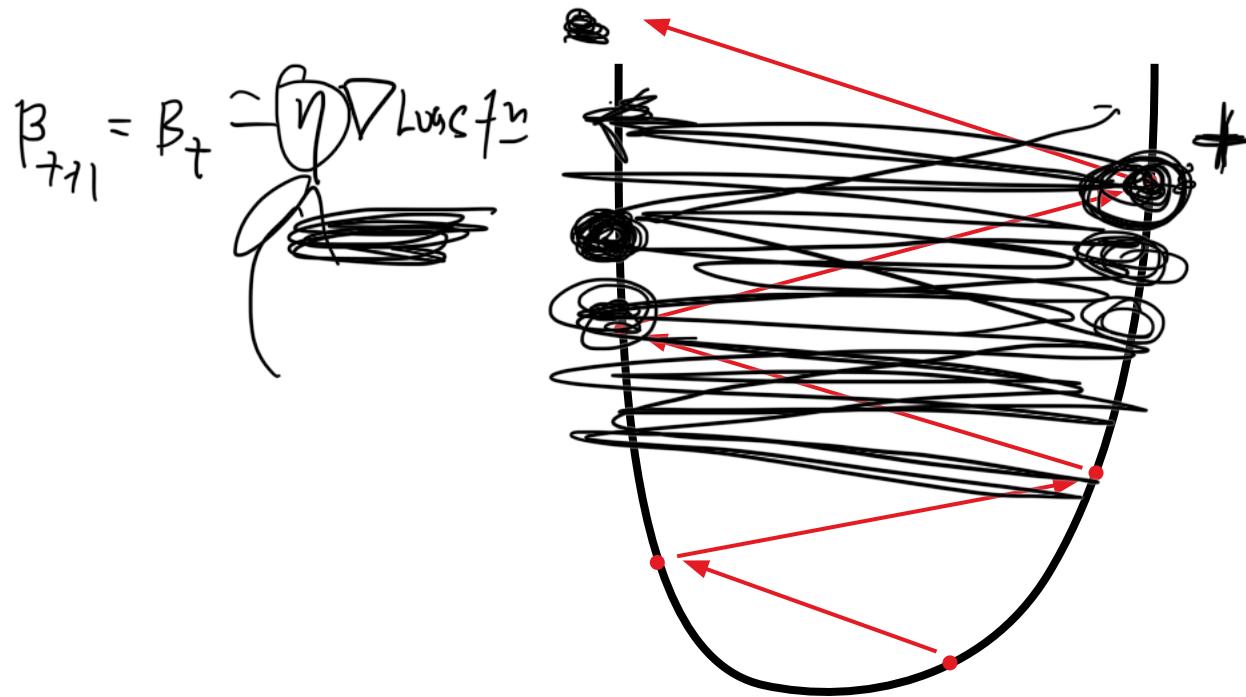


# Convexity

Thankfully, all\* loss functions in our course are convex!

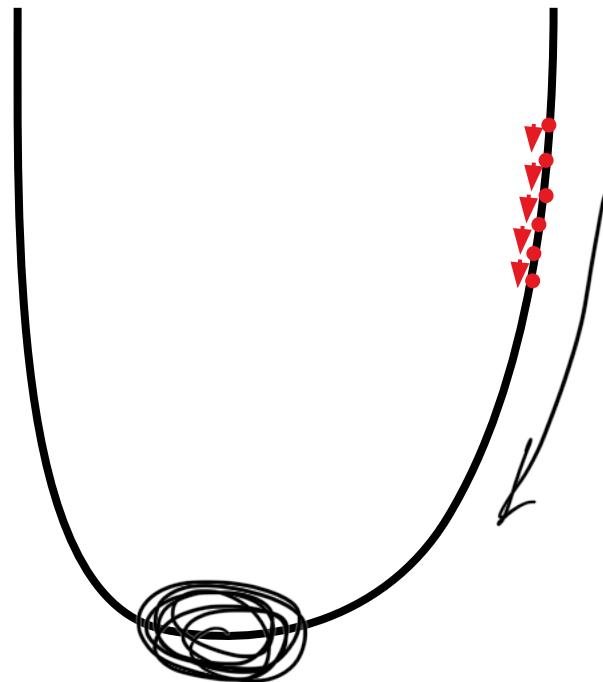
\*except the loss functions of neural networks. Always the exception...

## Problem #2: $\alpha$ too large

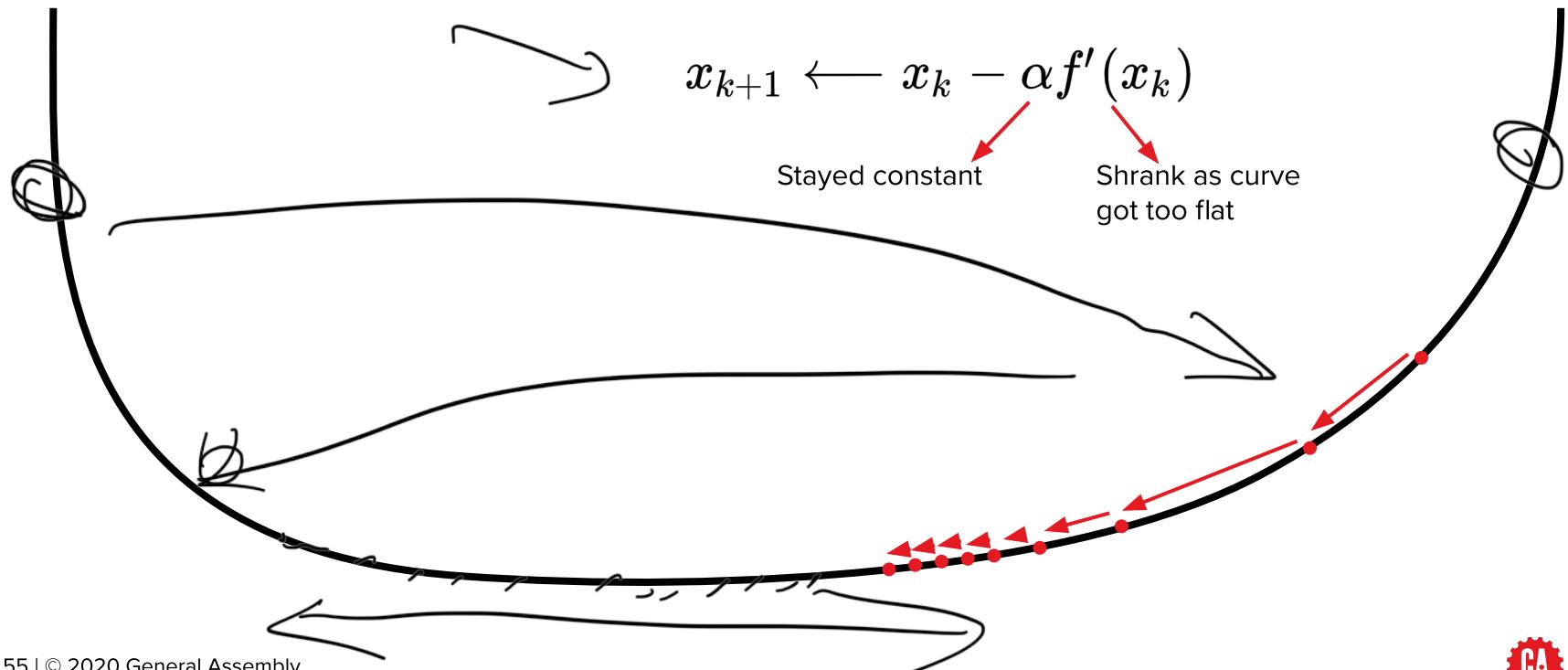


## Problem #3: $\alpha$ too small

convergenu



## Problem #4: Curve is too flat



—  
So then, how do we  
choose  $\alpha$ ?

## Choice of $\alpha$

---

Choosing the best value for our step size  $\alpha$  is probably the most difficult problem to solve. Here are a few ways:



## Method #1: Manually

$\eta$

We can choose  ~~$\alpha$~~  manually, and then using guess-and-check to iterate the process until we get an  $\alpha$  that converges.

It might not sound scientific, but it's what we'll do today. It's also what we'll do for neural networks, since other methods will save us from our folly.



## Method #2: Let $\alpha$ shrink over time

This is sensible, since we don't want  $\alpha$  being too large when it's time to converge. The following choices are common:

$$\alpha_k = \frac{c}{k} \text{ or } \frac{c}{\sqrt{k}}$$

While it's a good try, this rarely works in practice and adjusting c just amounts to doing guess-and-check again.

Adam

RMSPWD

## Method #3: Use a fancier algorithm

$$\beta_{t+1} = \beta_t - \alpha \nabla$$

The version of gradient descent we've learned today is overly simple for today's world, but it still gets to the gist of what we're trying to learn. In practice, more advanced methods are used, for example:

- Iteratively Reweighted Least Squares (IRLS)
- Adaptive Momentum (Adam) ←
- RMS Propagation (RMSProp) ←

## Method #3: Use a fancier algorithm

The version of gradient descent we've learned today is overly simple for today's world, but it still gets to the gist of what we're trying to learn. In practice, more advanced methods are used, for example:

- Iteratively Reweighted Least Squares (IRLS) ← **GLMs use this!**
- Adaptive Momentum (Adam) ←
- RMS Propagation (RMSProp) ← **Neural networks often use these!**

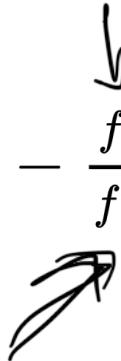


## Method #4: Second-order methods

On that note, what we've been doing is actually **first-order gradient descent** since it only uses the first derivative. Second-order gradient descent methods exist, the most common of which is known as **Newton's method** (yes, that Newton!) 

# Newton's Method

Instead of using a step size  $\alpha$ , we'll use:

$$x_{k+1} \leftarrow x_k - \frac{f'(x_k)}{f''(x_k)}$$


# Newton's Method

Instead of using a step size  $\alpha$ , we'll use:

$$x_{k+1} \leftarrow x_k - \frac{f'(x_k)}{f''(x_k)}$$

Not only does this fix all of our problems, but under most conditions, it's *guaranteed to work, and to work faster than first order methods!*

... so what's the catch?

# Newton's Method

$$x_{k+1} \leftarrow x_k - \frac{f'(x_k)}{f''(x_k)}$$

Turns out, second derivatives are **really** hard to calculate. Everything we've done today looks simple, but most problems are way harder to write out. Derivatives need to be estimated by computers - a whole other field of study!

Furthermore, most optimization problems work with vector-valued xs, which makes this orders of magnitude more complicated!

# Vector-valued Newton's method

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

# Vector-valued Newton's method

$$\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k - \underbrace{\left[ \nabla^2 f(\mathbf{x}_k) \right]^{-1}}_{\text{matrix}} \underbrace{\nabla f(\mathbf{x}_k)}_{\text{vector}}$$

Fancy symbol for  
“vector derivative”

That inverse is the tricky part. It's almost always computationally infeasible to find it. You're better off with method #3 at that point!

# Food for thought during our break

There are always students who don't believe me that there are simple one-variable problems that can't be optimized using calculus. If you are in that boat, convince yourself with this exercise! I'll even give you the derivative. We'll solve it with Python after the break.

$$f(x) = -\frac{\log x}{1+x}$$

$$f'(x) = -\frac{1/x+1-\log x}{(1+x)^2}$$