



Gradient Descent



Gradient Descent is one method for solving an optimization problem.

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

$$\min f(x)$$

we already know how to find minima/maxima.



Why We Need Gradient Descent

From calculus, we know that the minimum of f must lie at a point where its derivative vanishes, i.e. $\frac{\partial f(\theta^*)}{\partial \theta} = 0$.

- Sometimes, we can solve this equation analytically for θ .
- Mostly, we are not so lucky and must resort to iterative methods.

Recall the Gradient:

$$\nabla_{\theta} f = \left(\frac{\partial f(\theta)}{\partial \theta_1}, \frac{\partial f(\theta)}{\partial \theta_2}, \dots, \frac{\partial f(\theta)}{\partial \theta_n} \right)$$





$$f(x) = x^2 \cdot \log_e(x) - x$$

$$f'(x) = 2x \cdot \log_e(x) + x - 1 = 0 \Rightarrow x = \dots$$

$$\Rightarrow 2x \ln x + x = 1$$

x =

Not able to find
closed form.



Closed Form Solutions are not always available



- Often it is not possible to simply solve $\nabla_X f(X) = 0$
 - The function to minimize/maximize may have an intractable form
- In these situations, iterative solutions are used
 - Begin with a “guess” for the optimal X and refine it iteratively until the correct value is obtained

ES



Solving the equation $f'(x) = 0$ for x provides an *analytical* solution for a critical point. Unfortunately, it is not always possible to compute such analytical solutions in closed form. It is often difficult to exactly solve the equation $f'(x) = 0$ because this derivative might itself be a complex function of x . In other words, a *closed form solution* (like the example above) typically does not exist. For example, consider the following function that needs to be minimized:

$$f(x) = x^2 \cdot \log_e(x) - x \quad (4.5)$$

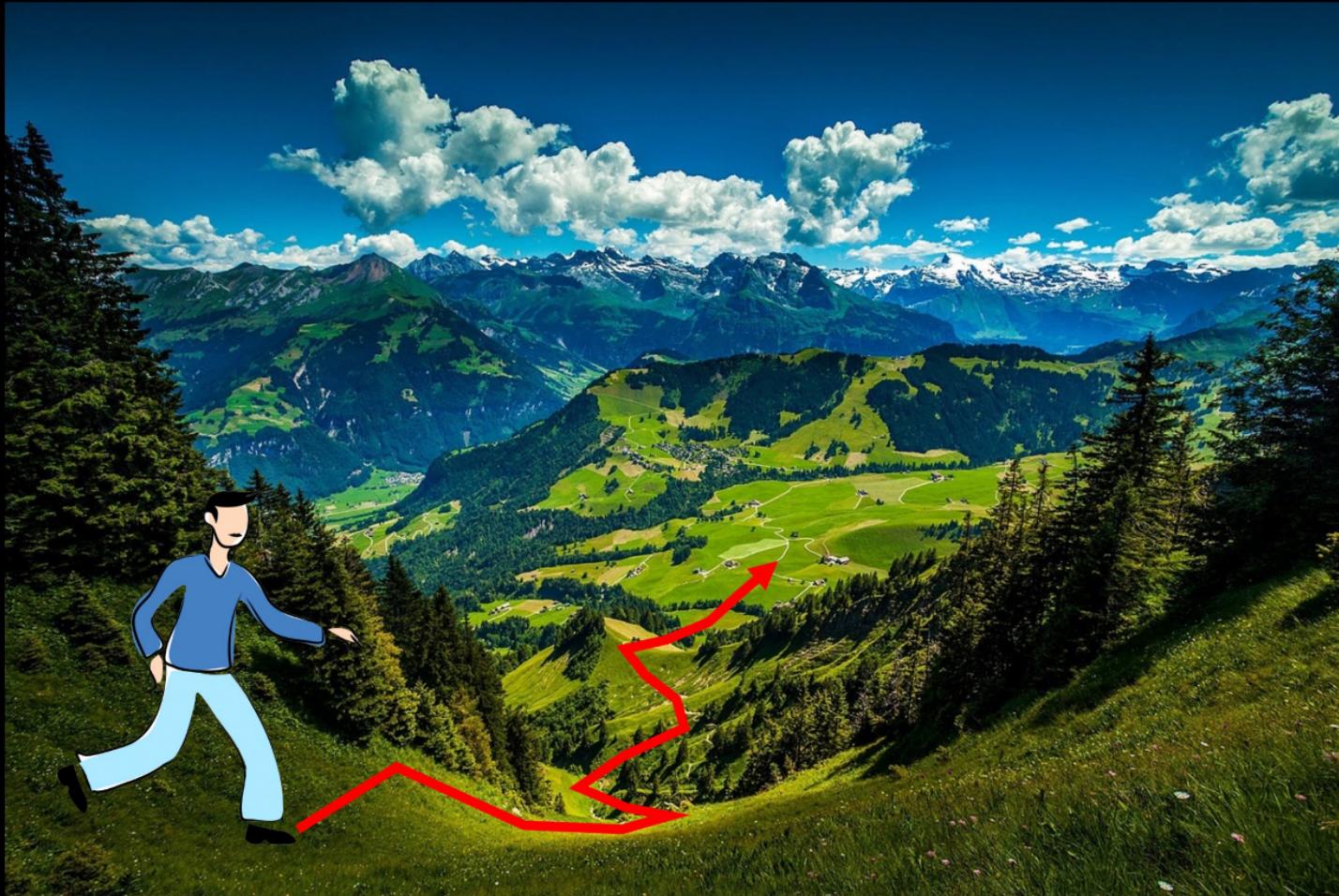
Setting the first derivative of this function to 0 yields the following condition:

$$f'(x) = 2x \cdot \log_e(x) + x - 1 = 0$$

This equation is somewhat hard to solve, although iterative methods exist for solving it. By trial and error, one might get lucky and find out that $x = 1$ is indeed a solution to the first-order optimality condition because it satisfies $f'(1) = 2 \log_e(1) + 1 - 1 = 0$. Furthermore, the second derivative $f''(x)$ can be shown to be positive at $x = 1$, and therefore this point is at least a local minimum. However, solving an equation like this numerically causes all types of numerical and computational challenges; these types of challenges increase when we move from univariate optimization to multivariate optimization.



Calculus



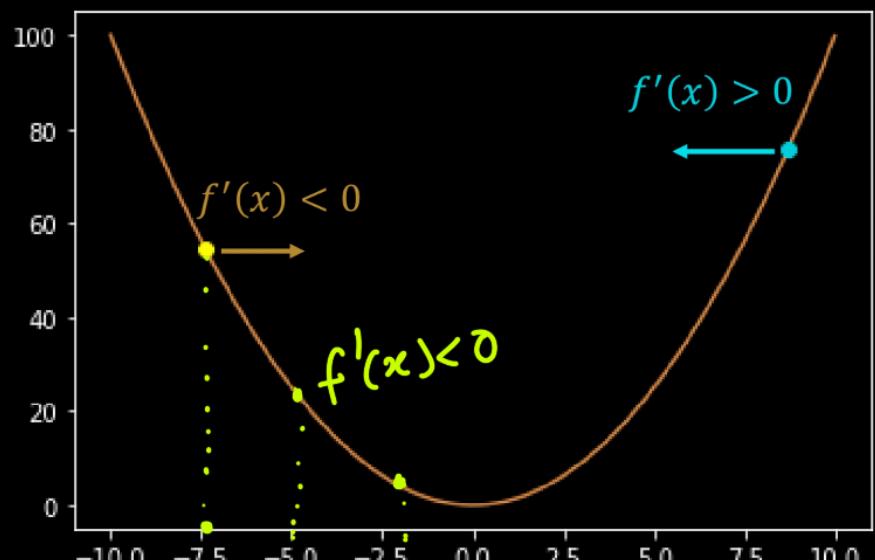


1D Gradient Descent

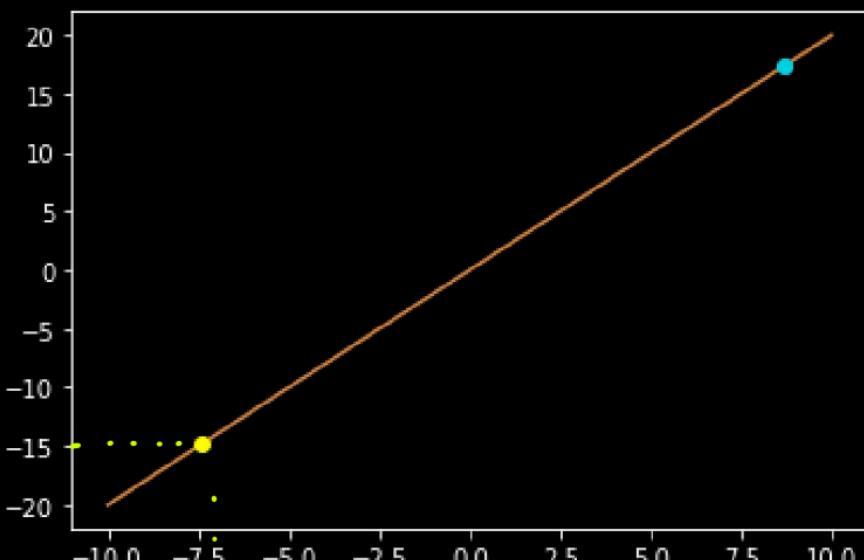


1D Gradient Descent

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



$$f'(x) = 2x$$



-7.5 ↗ -5.0 ↗ -2 ↗ -1 ↗ 0



1D Gradient Descent

$$x_1 = -7.5$$

$$x_{\text{new}} = x_1 + 2.5$$

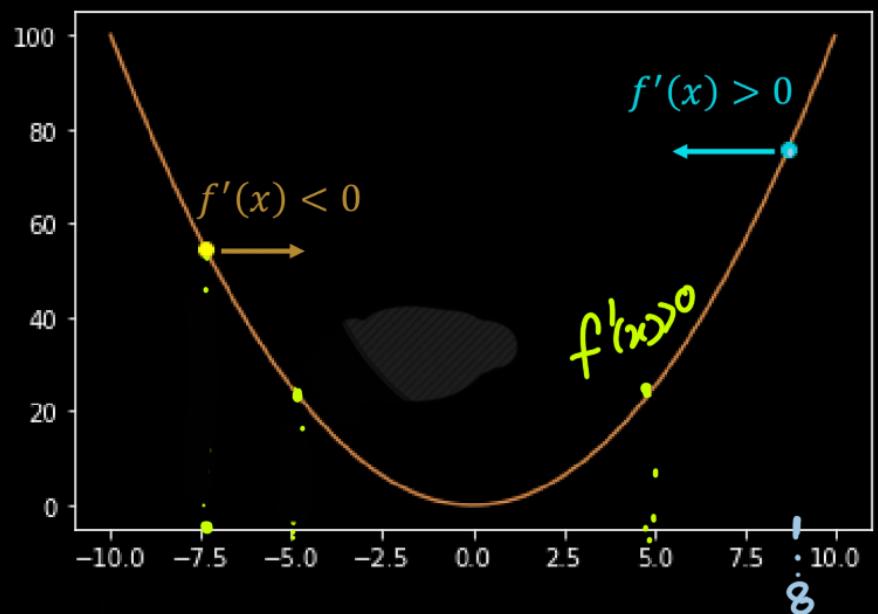
$$= -7.5 + 2.5$$

$$= -5.0$$

$$x_1 = 8$$

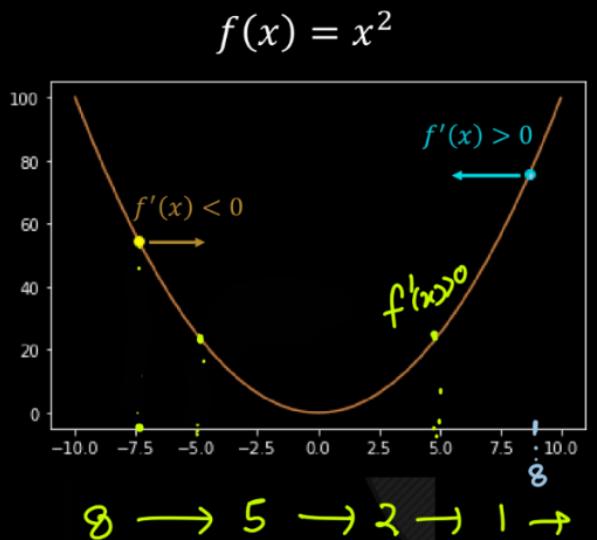
$$x_{\text{new}} = x_1 - 2.5$$

$$= 5.5$$



8 → 5 → 2 → 1 → 0

1D Gradient Descent



$$x_1 = 8$$

$$x_{\text{new}} = x_1 + 2.5$$

$$= 8 + 2.5$$

$$= 5.0$$

$$x_1 = 8$$

$$\begin{aligned} x_{\text{new}} &= x_1 - 2.5 \\ &= 5.5 \end{aligned}$$

step size or learning rate

$$x_{\text{new}} = x_1 - \alpha f'(x)$$

$$x_{\text{new}} = x_1 - \frac{1}{10} (-15)$$

$$= 8 + 1.5$$

$\underline{-f'(x)}$ ← direction of steepest decrease

$$\underline{x_{\text{new}} = x_1 - \alpha f'(x)}$$

- in which direction we should go
 - (to decrease the function)
- How much to go in that direction.

- in which direction we should go
 - (to decrease the function) - $f'(x)$
- How much to go in that direction.
 α



Gradient Descent (GD) (idea)



1. Start with a random value of w (e.g. $w = 12$)
2. Compute the gradient (derivative) of $L(w)$ at point $w = 12$. (e.g. $dL/dw = 6$)
3. Recompute w as:
$$w = w - \lambda * (dL / dw)$$



Gradient Descent

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

Step size: .8

$$x^{(0)} = -4$$

$$x_1 = x_0 - \alpha f'(x)$$

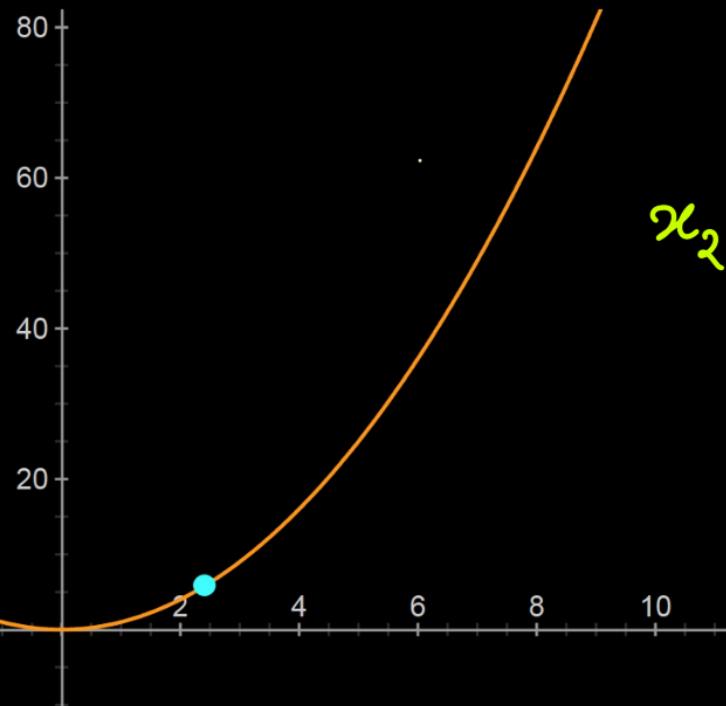
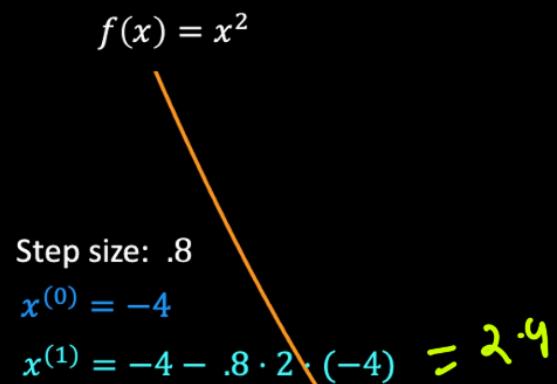
$$-4 - 0.8(-8)$$



$$= -4 + 0.8(-8) = \underline{\underline{2.4}}$$



Gradient Descent



$$x_{\text{new}} = x_{\text{old}} - \alpha f'_x(x)$$

$$x_2 = x_1 - 0.8 (2 \times 2.4)$$

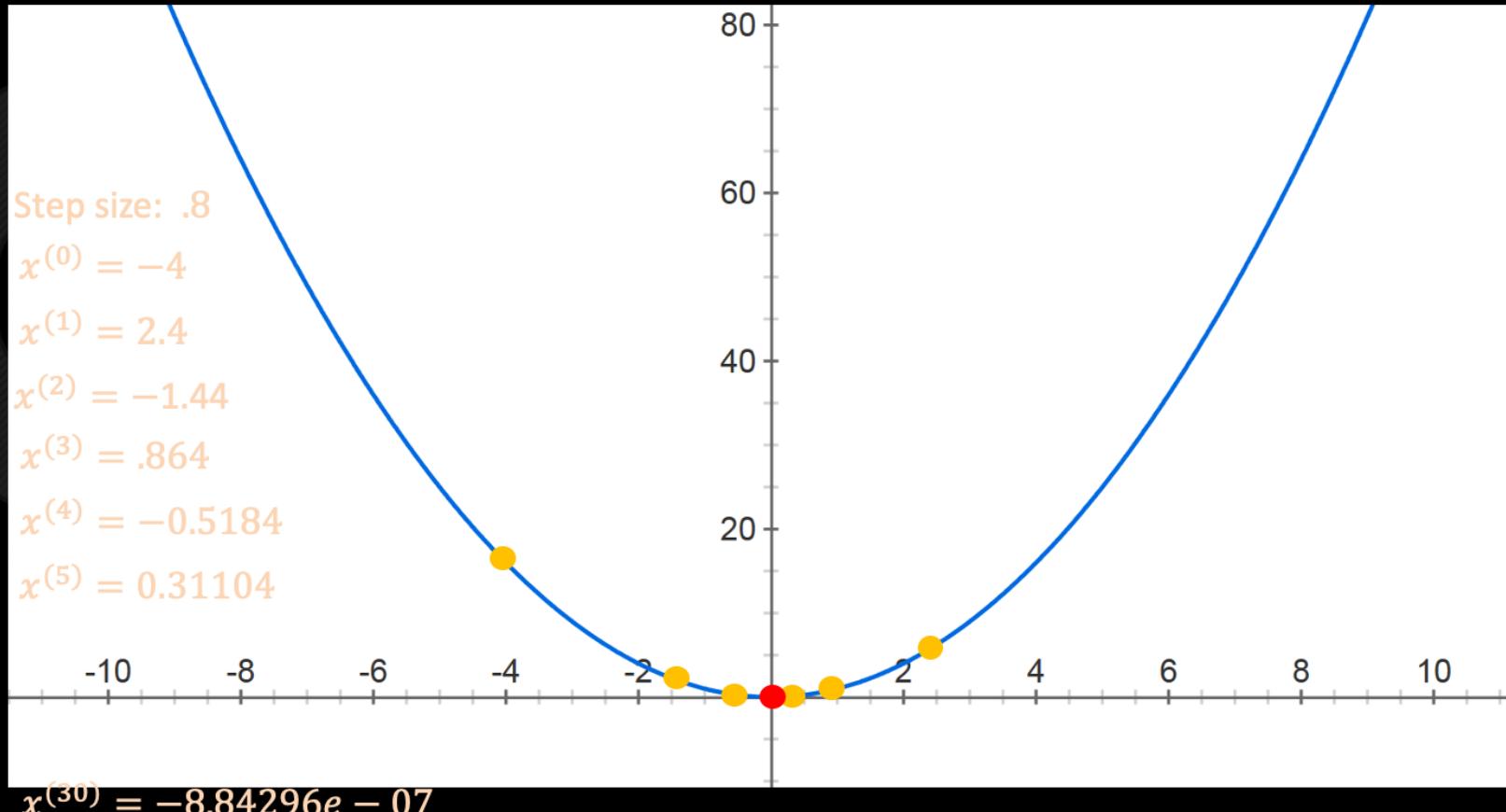
$$= 2.4 - 3.84$$

$$= -1.44$$

=====

Gradient Descent

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





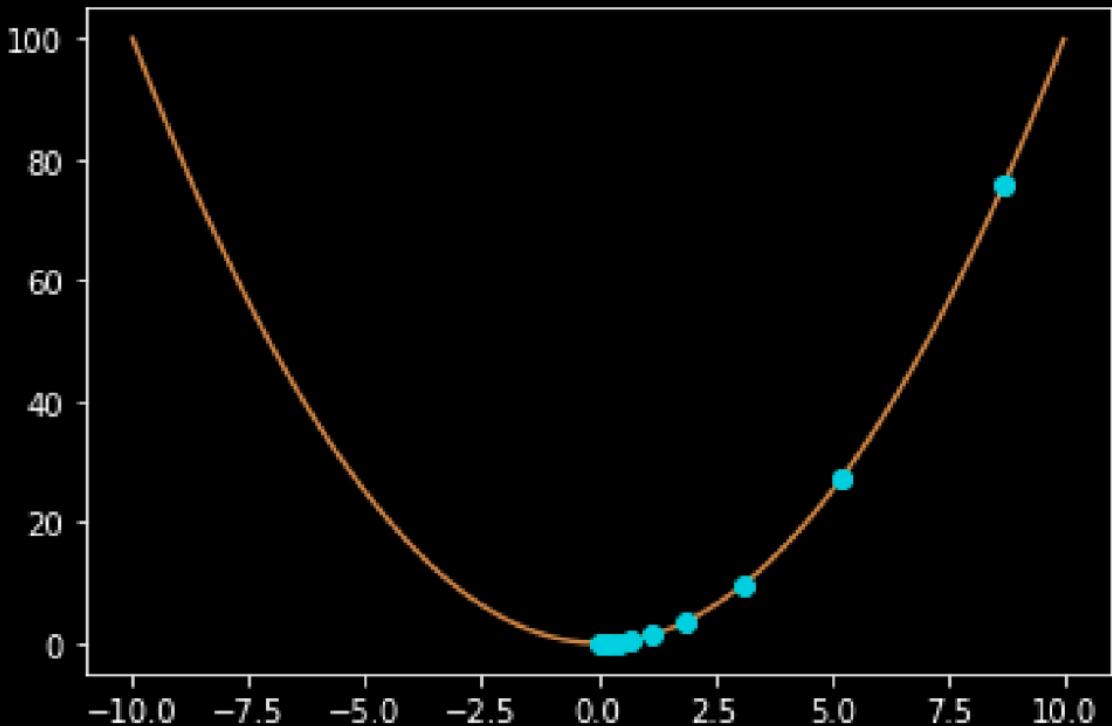
Calculus

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

$$f'(x) = 2x$$

$$x_0 = 8.7, \alpha = 0.2$$

$$x \leftarrow x - \alpha f'(x)$$





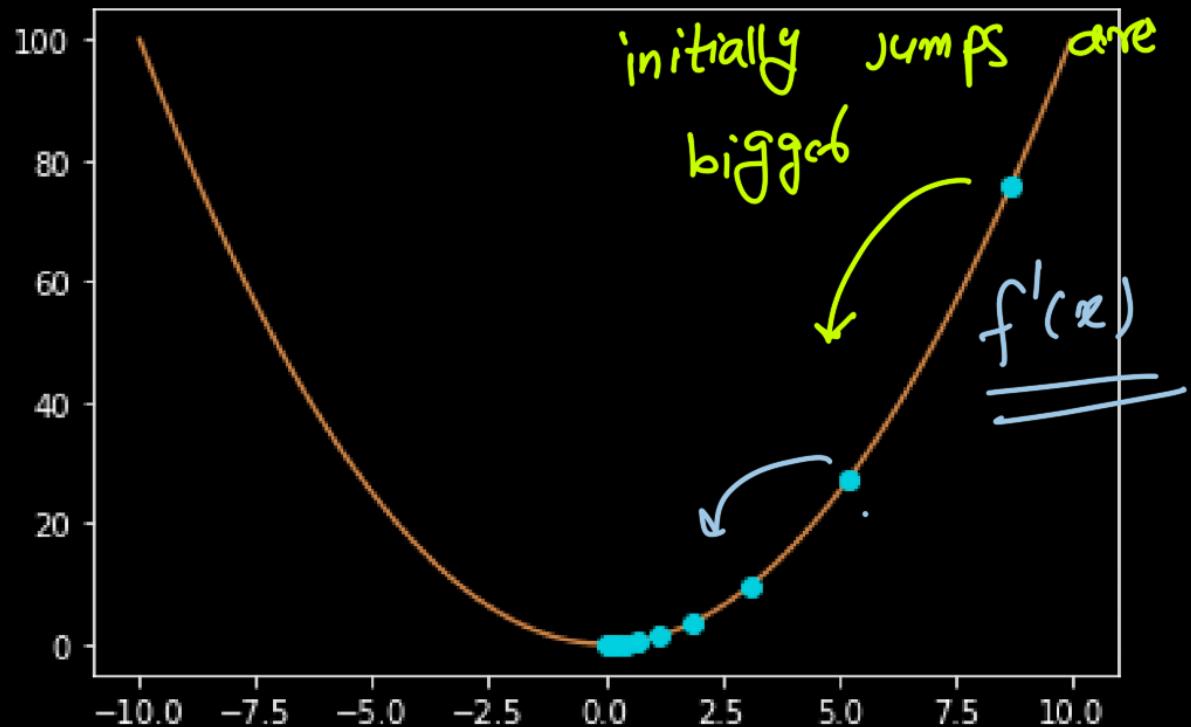
Calculus

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

$$f'(x) = 2x$$

$$x_0 = 8.7, \alpha = 0.2$$

$x \leftarrow x - \alpha f'(x)$





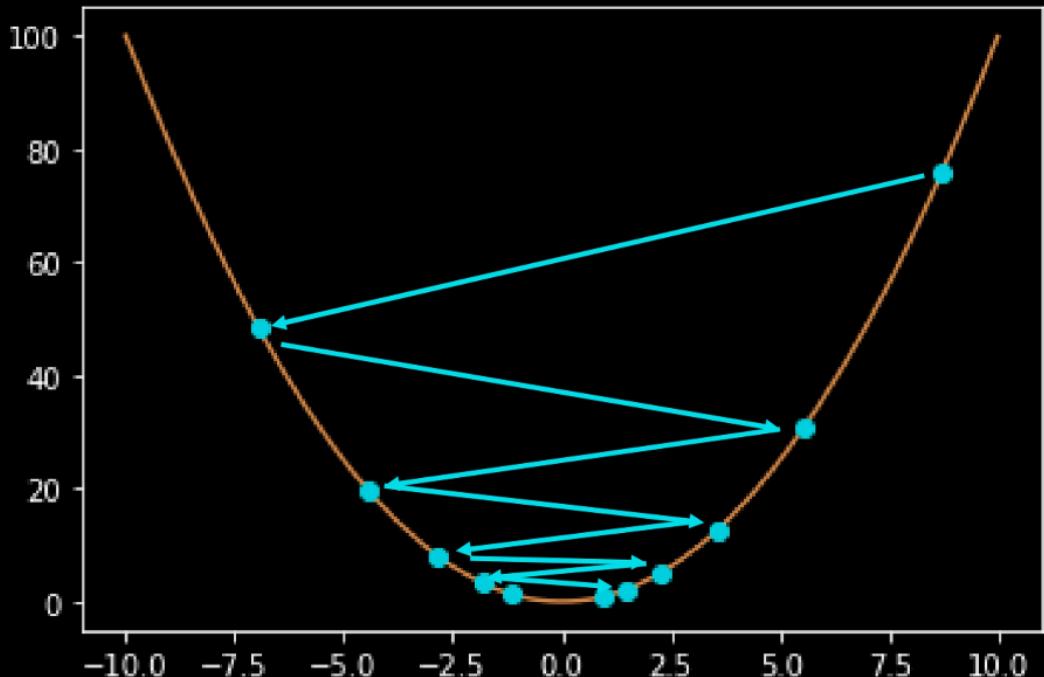
Calculus

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

$$f'(x) = 2x$$

$$x_0 = 8.7, \alpha = 0.9$$

$$x \leftarrow x - \alpha f'(x)$$





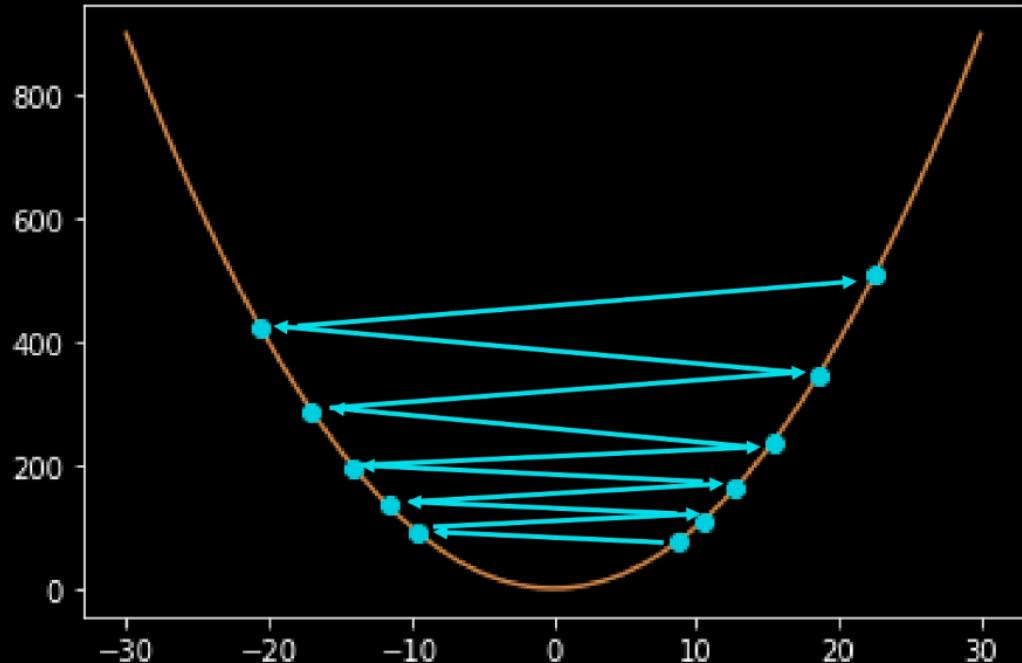
Calculus

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

$$f'(x) = 2x$$

$$x_0 = 8.7, \alpha = 1.05$$

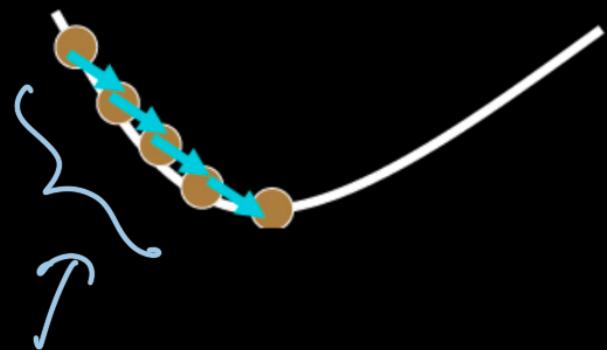
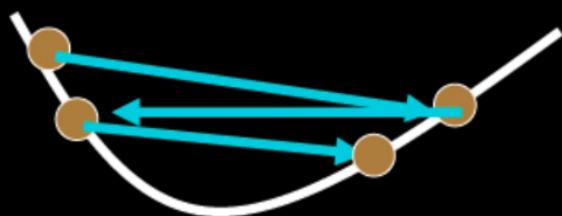
$$x \leftarrow x - \alpha f'(x)$$



Gradient descent may not even converge.



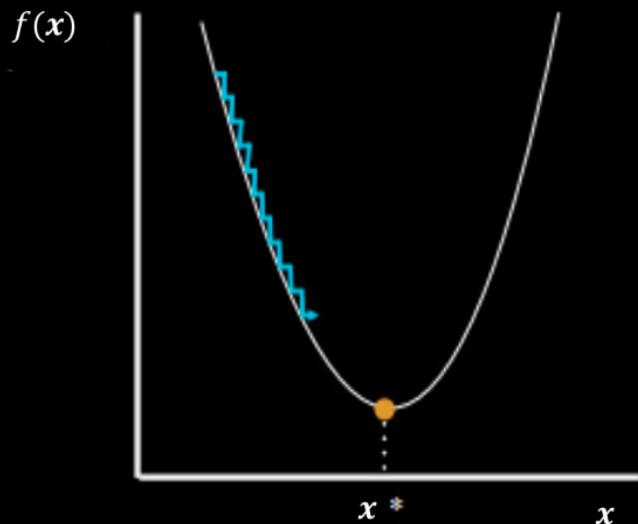
How to set step size?



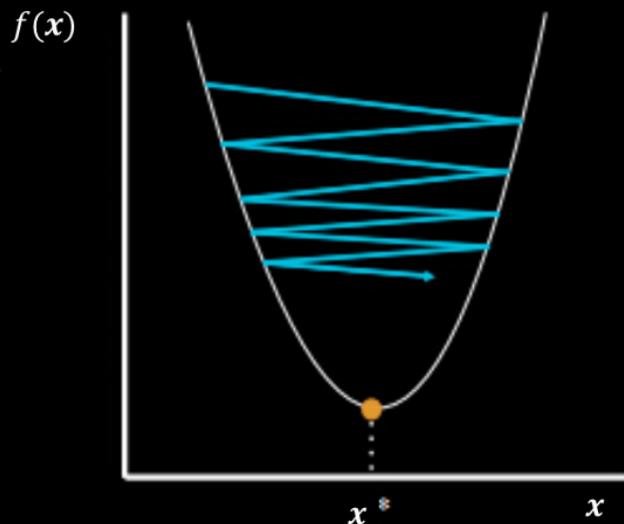
Slow walk



Choosing good step size matters!



Too small: converge
very slowly



Too big: overshoot and
even diverge

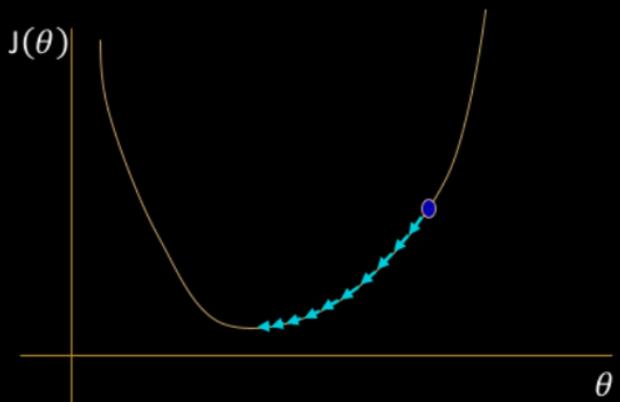
S





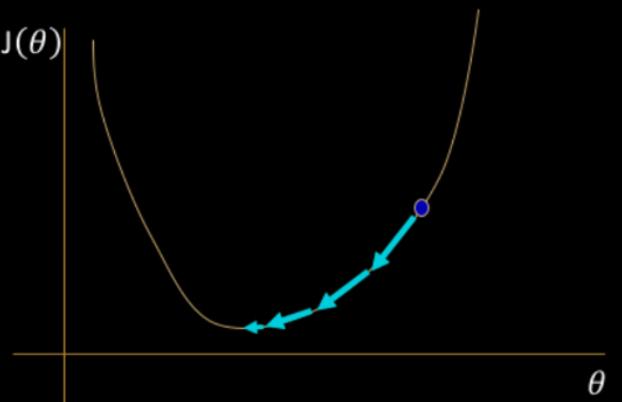
$$\Rightarrow \alpha = 0.1, 0.01, 0.2$$

Too low



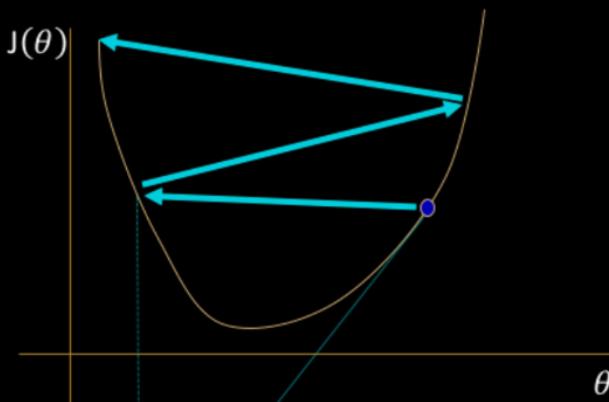
A small learning rate requires many updates before reaching the minimum point

Just right

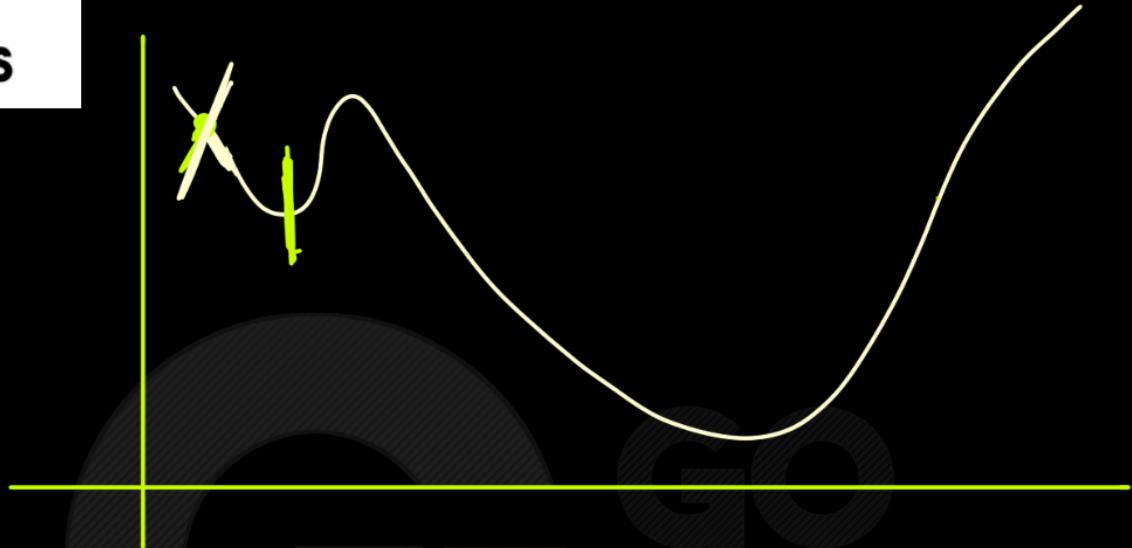


The optimal learning rate swiftly reaches the minimum point

Too high



Too large of a learning rate causes drastic updates which lead to divergent behaviors



algorithm may find local minima, not the
global minima.

↳ crucial factor on which it depends ⇒
Starting point.

The Approach of Gradient Descent



- Iterative solution:
 - Start at some point
 - Find direction in which to shift this point to decrease error
 - This can be found from the derivative of the function
 - A negative derivative → moving right decreases error
 - A positive derivative → moving left decreases error
 - Shift point in this direction



Gradient Descent

1. Initialize the parameters w to some guess
(usually all zeros, or random values)

2. Update the parameters:
 $w = w - \eta \nabla L(w)$

3. Update the learning rate η

4. Repeat steps 2-3 until $\nabla L(w)$ is close to zero.

stopping criteria

$$x_{\text{new}} = x_{\text{old}} - \eta f'(x)$$

ES

of number of iterations reaches 30



2D Gradient Descent

or multivariate



Gradient Descent

The algorithm:

$$\underline{x}^{(t+1)} = \underline{x}^{(t)} - \alpha \nabla f(\underline{x}^{(t)}), \quad t = 0, 1, 2, \dots,$$

where α is called the **step size**.

$$x_{\text{new}} = x_{\text{old}} - \cancel{\alpha \nabla f(x)}$$

$$\underline{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$x_i^{\text{new}} = x_i^{\text{old}} - \cancel{\alpha \frac{\partial f}{\partial x_i}}$

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha \nabla f(\mathbf{x}^{(t)})$$

$$\begin{bmatrix} x_1^{\text{new}} \\ \vdots \\ x_i^{\text{new}} \end{bmatrix} = \begin{bmatrix} x_1^{\text{old}} \\ \vdots \\ x_i^{\text{old}} \end{bmatrix} - \lambda \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_i} \end{bmatrix}$$



- Gradient descent is an **iterative algorithm**, which means we apply an update repeatedly until some criterion is met.
- We **initialize** the weights to something reasonable (e.g. all zeros) and repeatedly adjust them in the **direction of steepest descent**.



CLASSES

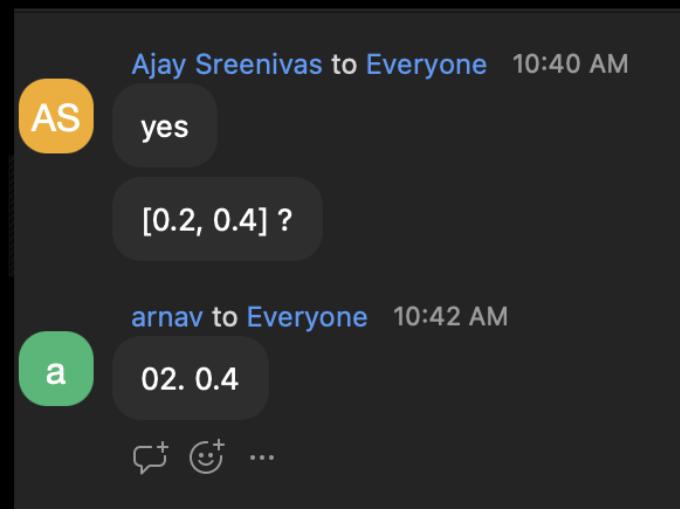


Question:

(2 points) Consider applying gradient descent with step size $\alpha = 0.1$ to find the \mathbf{x} that minimizes the function $f(\mathbf{x}) = f((x^{(1)}, x^{(2)})) = (x^{(1)} - 1)^2 + (x^{(2)} - 2)^2$ starting from $\mathbf{x}_0 = (0, 0)$. Find the value \mathbf{x}_1 after one iteration.

$$f(\mathbf{x}) = (x_1 - 1)^2 + (x_2 - 2)^2$$

$$\mathbf{x}_0 = (0, 0) \quad \alpha = 0.1$$



Question:

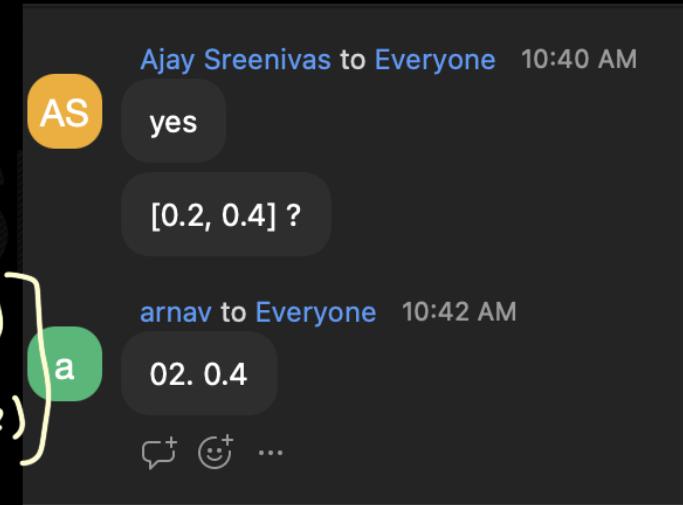
(2 points) Consider applying gradient descent with step size $\alpha = 0.1$ to find the \mathbf{x} that minimizes the function $f(\mathbf{x}) = f((x^{(1)}, x^{(2)})) = (x^{(1)} - 1)^2 + (x^{(2)} - 2)^2$ starting from $\mathbf{x}_0 = (0, 0)$. Find the value \mathbf{x}_1 after one iteration.

$$f(\mathbf{x}) = (x_1 - 1)^2 + (x_2 - 2)^2$$

$$\mathbf{x}_0 = (0, 0) \quad \alpha = 0.1$$

$$\mathbf{x}_1 = \mathbf{x}_0 - \alpha \nabla f(\mathbf{x})$$

$$\begin{aligned}\mathbf{x}_1 &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \begin{bmatrix} -2 \\ -4 \end{bmatrix} \\ &= \begin{bmatrix} 0.2 \\ 0.4 \end{bmatrix}\end{aligned}$$



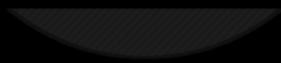


(2 points) Consider applying gradient descent with step size $\alpha = 0.1$ to find the \mathbf{x} that minimizes the function $f(\mathbf{x}) = f((x^{(1)}, x^{(2)})) = (x^{(1)} - 1)^2 + (x^{(2)} - 2)^2$ starting from $\mathbf{x}_0 = (0, 0)$. Find the value \mathbf{x}_1 after one iteration.

ANSWER:

$\nabla f(\mathbf{x}) = (2(x^{(1)} - 1), 2(x^{(2)} - 2))$, which is $(-2, -4)$ at $\mathbf{x} = (0, 0)$.

Move to $\mathbf{x}_1 = \mathbf{x}_0 - \alpha \nabla f(\mathbf{x}_0) = (0, 0) - (-0.2, -0.4) = (0.2, 0.4)$.





Question:

(8 points) Mark each statement as true (T) or false (F).

- Gradient descent can fail to converge on a convex function if step size α is such that we get stuck in a cycle, oscillating between two or several values.
 - For a non-convex function, gradient descent can fail to converge by descending without bound.
 - Gradient descent can fail to converge on a convex function if it gets stuck in a local minimum.
 - Gradient descent can fail to converge on a convex function if the step size $\alpha > 0$ is too small.
-



Question:

(8 points) Mark each statement as true (T) or false (F).

- T** _____ Gradient descent can fail to converge on a convex function if step size α is such that we get stuck in a cycle, oscillating between two or several values.
- _____ For a non-convex function, gradient descent can fail to converge by descending without bound.
- F** _____ Gradient descent can fail to converge on a convex function if it gets stuck in a local minimum.
- F** _____ Gradient descent can fail to converge on a convex function if the step size $\alpha > 0$ is too small.



| (8 points) Mark each statement as true (T) or false (F).

- Gradient descent can fail to converge on a convex function if step size α is such that we get stuck in a cycle, oscillating between two or several values. **ANSWER: T**
 - For a non-convex function, gradient descent can fail to converge by descending without bound. **ANSWER: T**
 - Gradient descent can fail to converge on a convex function if it gets stuck in a local minimum. **ANSWER: F**
 - Gradient descent can fail to converge on a convex function if the step size $\alpha > 0$ is too small. **ANSWER: F**
-



Question:

H.w.
==

Consider applying gradient descent with step size $\alpha = 0.5$ to find the \mathbf{x} that minimizes the function $f(\mathbf{x}) = f((x^{(1)}, x^{(2)})) = (x^{(1)} - 3)^2 + (x^{(2)} - 4)^2$ starting from $\mathbf{x}_0 = (1, 2)$. Find the value \mathbf{x}_1 after one iteration.





Calculus

Consider applying gradient descent with step size $\alpha = 0.5$ to find the \mathbf{x} that minimizes the function $f(\mathbf{x}) = f((x^{(1)}, x^{(2)})) = (x^{(1)} - 3)^2 + (x^{(2)} - 4)^2$ starting from $\mathbf{x}_0 = (1, 2)$. Find the value \mathbf{x}_1 after one iteration.

ANSWER:

$$\nabla f(\mathbf{x}) = (2(x^{(1)} - 3), 2(x^{(2)} - 4)), \text{ which is } (-4, -4) \text{ at } \mathbf{x} = (1, 2).$$

$$\text{Move to } \mathbf{x}_1 = \mathbf{x}_0 - \alpha \nabla f(\mathbf{x}_0) = (1, 2) - 0.5(-4, -4) = (3, 4).$$



Question:

[3 Pts] Suppose we run gradient descent with a fixed learning rate of $\alpha = 0.1$ to minimize the 2D function $f(x, y) = 5 + x^2 + y^2 + 5xy$.

The gradient of this function is

$$\nabla_{x,y} f(x, y) = \begin{bmatrix} 2x + 5y \\ 2y + 5x \end{bmatrix} = \begin{bmatrix} 12 \\ 9 \end{bmatrix}$$

If our starting guess is $x^{(0)} = 1, y^{(0)} = 2$, what will be our next guess $x^{(1)}, y^{(1)}$?

$$x^{(1)} = \boxed{}$$

$$y^{(1)} = \boxed{} = \begin{bmatrix} -0.2 \\ 1.1 \end{bmatrix}$$



Calculus

$$x^{(1)} = \boxed{}$$

$$y^{(1)} = \boxed{}$$

Solution: The gradient is $= [2*1 + 5*2, 2*2 + 5*1] = [12, 9]$ so next guess is $[1, 2] - 0.1 * [12, 9] = -0.2, 1.1$





Question:

Consider the following loss function on vectors $w \in \mathbb{R}^4$:

$$L(w) = w_1^2 + 2w_2^2 + w_3^2 - 2w_3w_4 + w_4^2 + 2w_1 - 4w_2 + 4.$$

- (a) What is $\nabla L(w)$?
- (b) Suppose we use gradient descent to minimize this function, and that the current estimate is $w = (0, 0, 0, 0)$. If the step size is η , what is the next estimate?
- (c) What is the minimum value of $L(w)$?
- (d) Is there is a unique solution w at which this minimum is realized?

$$\nabla L(w) = (2w_1 + 2, 4w_2 - 4, 2w_3 - 2w_4, -2w_3 + 2w_4)$$



Question:

Consider the following loss function on vectors $w \in \mathbb{R}^4$:

$$L(w) = w_1^2 + 2w_2^2 + w_3^2 - 2w_3w_4 + w_4^2 + 2w_1 - 4w_2 + 4.$$

- (a) What is $\nabla L(w)$?
- (b) Suppose we use gradient descent to minimize this function, and that the current estimate is $w = (0, 0, 0, 0)$. If the step size is η , what is the next estimate?
- (c) What is the minimum value of $L(w)$?
- (d) Is there is a unique solution w at which this minimum is realized?

$$\nabla L(w) = (2w_1 + 2, 4w_2 - 4, 2w_3 - 2w_4, -2w_3 + 2w_4)$$

$$(2, -4, 0, 0) \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} - \eta \begin{pmatrix} 2 \\ -4 \\ 0 \\ 0 \end{pmatrix}$$



$$L(w) = w_1^2 + 2w_2^2 + w_3^2 - 2w_3w_4 + w_4^2 + 2w_1 - 4w_2 + 4$$

(a) The derivative is

$$\nabla L(w) = (2w_1 + 2, 4w_2 - 4, 2w_3 - 2w_4, -2w_3 + 2w_4)$$

(b) The derivative at $w = (0, 0, 0, 0)$ is $(2, -4, 0, 0)$. Thus the update at this point is:

$$w_{new} = w - \eta \nabla L(w) = (0, 0, 0, 0) - \eta(2, -4, 0, 0) = (-2\eta, 4\eta, 0, 0).$$

$$\underline{\eta = \frac{1}{2}}$$

(c) To find the minimum value of $L(w)$, we will equate $\nabla L(w)$ to zero:

- $2w_1 + 2 = 0 \implies w_1 = -1$
- $4w_2 - 4 = 0 \implies w_2 = 1$
- $2w_3 - 2w_4 = 0 \implies w_3 = w_4$

$$0, 0, 0, 0$$



$$\underline{-2\eta, 4\eta, 0, 0}$$

$$-1, 1, 0, 0$$

(d) No, there is not a unique solution.



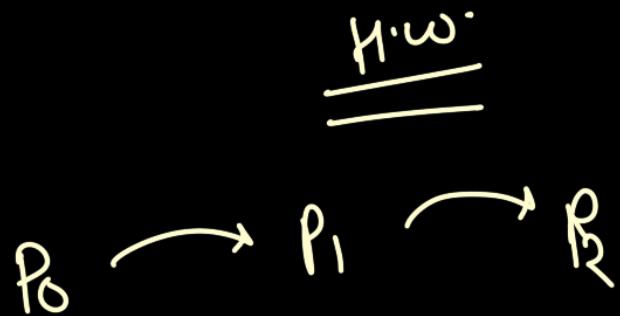
Question:

Let

$$f(x, y) = (x - 1)^2(y - 3)^2$$

$$P_0 = (x_0, y_0) = \underline{(4, 5)}$$

$$\eta = 0.05$$



Using the gradient descent algorithm find $P_2 = (x_2, y_2)$ and then calculate $f(P_2)$ to at least the nearest thousandth (at least 3 decimal places).

Answer

$$\nabla f = \left\langle \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right\rangle = \left\langle 2(x-1)(y-3)^2, 2(x-1)^2(y-3) \right\rangle$$

$$P_0 = (4, 5)$$

$$P_1 = P_0 - \eta \cdot \nabla f(P_0)$$

$$= (4, 5) - (0.05) \nabla f(4, 5)$$

$$= (4, 5) - (0.05) \left\langle 2(4-1)(5-3)^2, 2(4-1)^2(5-3) \right\rangle$$

$$= (4, 5) - (0.05) \left\langle 24, 36 \right\rangle$$

$$\underbrace{P_1}_{= (2.8, 3.2)} = (4, 5) - (0.05) \left\langle 24, 36 \right\rangle$$

$$P_2 = P_1 - \eta \cdot \nabla f(P_1)$$

$$= (2.8, 3.2) - (0.05) \nabla f(2.8, 3.2)$$

$$= (2.8, 3.2) - (0.05) \left\langle 2(2.8-1)(3.2-3)^2, 2(2.8-1)^2(3.2-3) \right\rangle$$

$$\underbrace{P_2}_{= (2.7928, 3.152)} = (2.8, 3.2) - (0.05) \left\langle 2(1.8)(-1)^2(3.2-3), 2(1.8)^2(-1) \right\rangle$$

$$f(P_2) = f(2.7928, 3.152)$$

$$= (2.7928-1)^2(3.152-3)^2$$

$$= 0.0588$$

answer

Given:

$$f(x, y) = (x-1)^2(y-3)^2$$

$$P_0 = \begin{pmatrix} x_0 \\ y_0 \end{pmatrix} = \begin{pmatrix} 4 \\ 5 \end{pmatrix}$$

$$\eta = 0.05$$

The gradient of $f(x, y)$ is:

$$\nabla f(x, y) = \begin{pmatrix} 2(x-1)(y-3)^2 \\ 2(x-1)^2(y-3) \end{pmatrix}$$

Calculate $\nabla f(4, 5)$:

$$\nabla f(4, 5) = \begin{pmatrix} 2 \cdot (4-1) \cdot (5-3)^2 \\ 2 \cdot (4-1)^2 \cdot (5-3) \end{pmatrix} = \begin{pmatrix} 24 \\ 36 \end{pmatrix}$$

Now, update P_1 :

$$P_1 = P_0 - \eta \cdot \nabla f(4, 5) = \begin{pmatrix} 4 \\ 5 \end{pmatrix} - 0.05 \cdot \begin{pmatrix} 24 \\ 36 \end{pmatrix} = \begin{pmatrix} 4 \\ 5 \end{pmatrix} - \begin{pmatrix} 1.2 \\ 1.8 \end{pmatrix} = \begin{pmatrix} 2.8 \\ 3.2 \end{pmatrix}$$

Calculate $\nabla f(2.8, 3.2)$:

$$\nabla f(2.8, 3.2) = \begin{pmatrix} 2 \cdot (2.8-1) \cdot (3.2-3)^2 \\ 2 \cdot (2.8-1)^2 \cdot (3.2-3) \end{pmatrix} = \begin{pmatrix} 0.072 \\ 0.648 \end{pmatrix}$$

Update P_2 :

$$P_2 = P_1 - \eta \cdot \nabla f(2.8, 3.2) = \begin{pmatrix} 2.8 \\ 3.2 \end{pmatrix} - 0.05 \cdot \begin{pmatrix} 0.072 \\ 0.648 \end{pmatrix} = \begin{pmatrix} 2.8 \\ 3.2 \end{pmatrix} - \begin{pmatrix} 0.0036 \\ 0.0324 \end{pmatrix} = \begin{pmatrix} 2.7964 \\ 3.1676 \end{pmatrix}$$

Finally, calculate $f(P_2)$:

$$f(P_2) = (2.7964-1)^2 \cdot (3.1676-3)^2$$



Question:

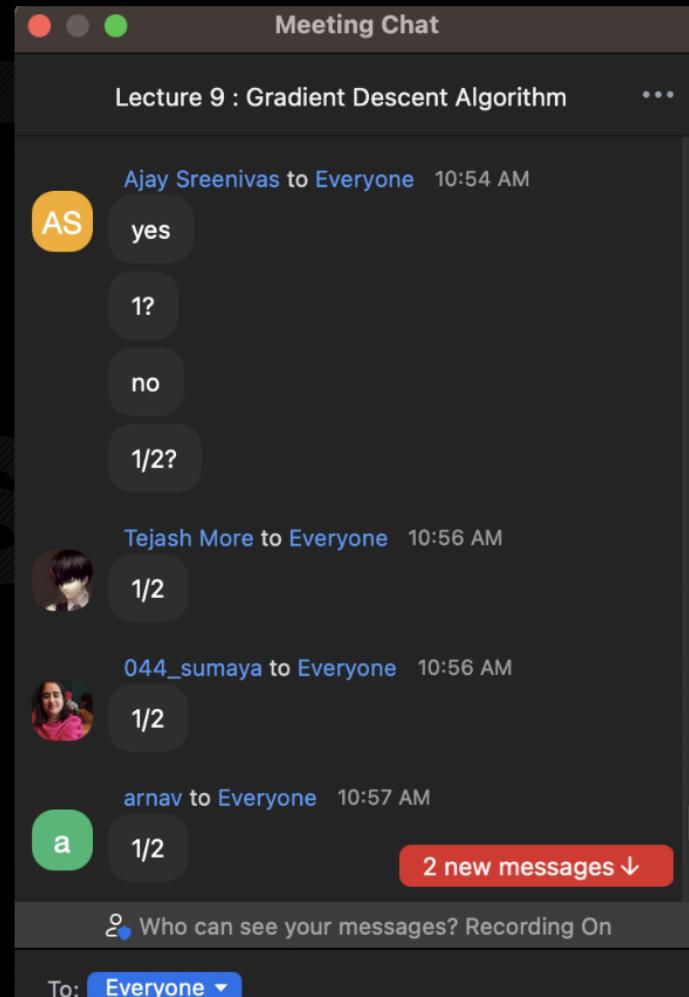
Consider the univariate function $f(x) = x^2$. This function has a unique minimum at $x^* = 0$. We're using gradient descent (GD) to find this minimum and at time t we arrive at the point $x_t = 2$. What is the step size that would bring us to x^* at time $t + 1$?



Consider the univariate function $f(x) = x^2$. This function has a unique minimum at $x^* = 0$. We're using gradient descent (GD) to find this minimum and at time t we arrive at the point $x_t = 2$. What is the step size that would bring us to x^* at time $t + 1$?

$$\begin{aligned}x_t &= 2 \\&\text{---} \\x_{t+1} &= 2 - \alpha f'(x) \\0 &= 2 - \alpha \cdot (4) \\&\Rightarrow \alpha = \frac{1}{2}\end{aligned}$$

$$f'(x) = 2x$$





Explanation: Using the definition of GD and the requirements in the problem statement we are looking for η such that:

$$x_* = x_t - \eta \nabla f(x_t) \iff 0 = 2 - \eta \nabla f(2) \iff \eta \cdot (2 \cdot 2) = 2 \iff \eta = \frac{1}{2}$$





Question:

(5 points) Consider the following loss function:

$$L(w) = \sum_{i=1}^n \log(1 + w^\top x_i)$$

$$\nabla L_w = \sum_{i=1}^n \nabla L_i$$

What is $\nabla L(w)$? Write down the update step for gradient descent.

$$L_i = \log(1 + w^\top x_i)$$

$$\nabla_w L_i = \frac{1}{1 + w^\top x_i} \cdot x_i$$

$$L(w) = \sum_{i=1}^n L_i$$

$$w_{\text{new}} = w_{\text{old}} - \eta \cdot \sum_{i=1}^n \nabla L_i$$



Calculus

(5 points) Consider the following loss function:

$$L(w) = \sum_{i=1}^n \log(1 + w^\top x_i)$$

What is $\nabla L(w)$? Write down the update step for gradient descent.

Let $L_i(w) = \log(1 + w^\top x_i)$. Then

$$\begin{aligned}\nabla L(w) &= \sum_{i=1}^n \nabla L_i(w) \\ &= \sum_{i=1}^n \frac{x_i}{1 + w^\top x_i}.\end{aligned}$$

The gradient descent update step for w is then

$$w_{t+1} = w_t - \eta_t \nabla L(w) = w_t - \eta_t \sum_{i=1}^n \frac{x_i}{1 + w_t^\top x_i},$$



where η_t is the learning rate at the t -th descent step.



Question:

MSQ

[3 Pts] The learning rate can *potentially* affect which of the following? Select all that apply. Assume nothing about the function being minimized other than that its gradient exists. You may assume the learning rate is positive.

- The speed at which we converge to a minimum.
- Whether gradient descent converges.
- The direction in which the step is taken.
- Whether gradient descent converges to a local minimum or a global minimum.

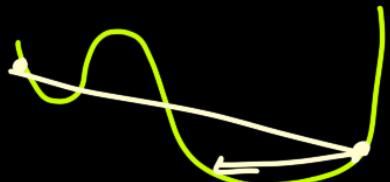


Question:

MSQ

[3 Pts] The learning rate can *potentially* affect which of the following? Select all that apply. Assume nothing about the function being minimized other than that its gradient exists. You may assume the learning rate is positive.

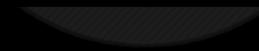
- The speed at which we converge to a minimum.
- Whether gradient descent converges.
- The direction in which the step is taken.
- Whether gradient descent converges to a local minimum or a global minimum.





[3 Pts] The learning rate can *potentially* affect which of the following? Select all that apply. Assume nothing about the function being minimized other than that its gradient exists. You may assume the learning rate is positive.

- The speed at which we converge to a minimum.**
- Whether gradient descent converges.**
- The direction in which the step is taken.
- Whether gradient descent converges to a local minimum or a global minimum.**





Question:

35. For the dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ and the loss function:

$$L(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \sin(w_0 + w_1 x_i))^2$$

H'w

answer each of the following questions.

Which of the following is the gradient of the loss function with respect to $w = [w_0, w_1]$

- 0
- $\frac{1}{n} \sum_{i=1}^n (y_i - \sin(w_0 + w_1 x_i))^2$
- $\frac{1}{n} \sum_{i=1}^n \sin(y_i - \cos(w_0 + w_1 x_i)) w_1$
- $-\frac{2}{n} \sum_{i=1}^n (y_i - \sin(w_0 + w_1 x_i)) \cos(w_0 + w_1 x_i) [1, x_i]$
- $[-\frac{2}{n} \sum_{i=1}^n \cos(w_0 + w_1 x_i), -\frac{2}{n} \sum_{i=1}^n \cos(w_0 + w_1 x_i) x_i]$



(a) [3 Pts] Which of the following is the gradient of the loss function with respect to $w = [w_0, w_1]$

- 0
- $\frac{1}{n} \sum_{i=1}^n (y_i - \sin(w_0 + w_1 x_i))^2$
- $\frac{1}{n} \sum_{i=1}^n \sin(y_i - \cos(w_0 + w_1 x_i) w_1)$
- $-\frac{2}{n} \sum_{i=1}^n (\sin(w_0 + w_1 x_i)) \cos(w_0 + w_1 x_i) [1, x_i]$
- $[-\frac{2}{n} \sum_{i=1}^n \cos(w_0 + w_1 x_i), -\frac{2}{n} \sum_{i=1}^n \cos(w_0 + w_1 x_i) x_i]$



Question:

Which of the following could affect whether or not gradient descent converges to the global minimum?

- (a) [1 Pt] Learning Rate

True False

- (b) [1 Pt] Initialization of parameters

True False

- (c) [1 Pt] The ordering of the data (ignoring issues related to numerical precision).

True False



Question:

Which of the following could affect whether or not gradient descent converges to the global minimum?

- (a) [1 Pt] Learning Rate

True False



learning rate can make it diverge

- (b) [1 Pt] Initialization of parameters

True False

crucial thing

- (c) [1 Pt] The ordering of the data (ignoring issues related to numerical precision).

True False





Calculus

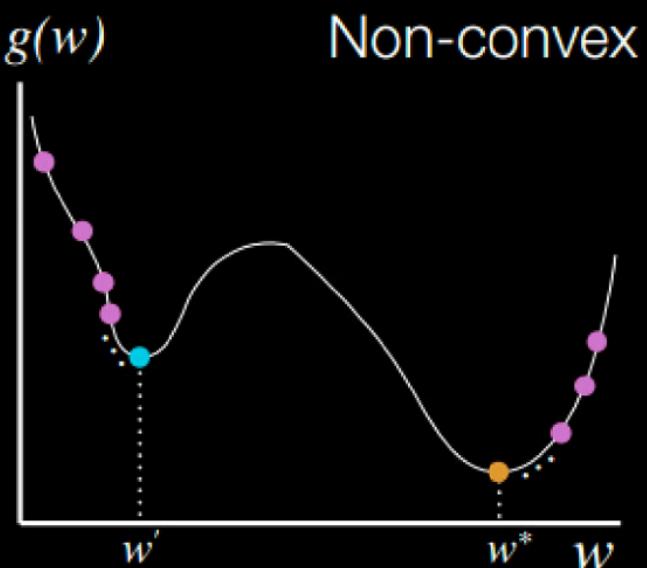
36.

- (a) [1 Pt] Learning Rate
 True False

- (b) [1 Pt] Initialization of parameters
 True False

- (c) [1 Pt] The ordering of the data (igno:
 True **False**

Solution: The initialization of parameters affects the starting parameters of gradient descent. If the loss function is not convex, the choice of starting parameters will affect whether the gradient descent algorithm converges to a local minimum or global minimum. The graph below shows a non-convex loss function with respect to parameter w .

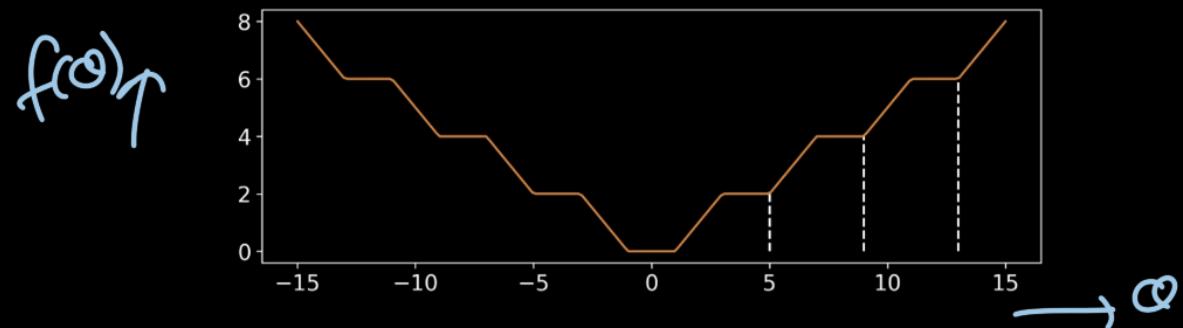


If we initialize w to be very small, then gradient descent will converge to a local minimum (given a reasonable learning rate). If we initialize w to be very large, then gradient descent will converge to a global minimum (given a reasonable learning rate).



Question:

13. Consider the following function of $f(\theta)$, which alternates between completely flat regions and regions of absolute slope equal to 1. **There is only one correct answer for each part.**



- (a) Assuming that θ starts in a flat region that is not a minimum and $\alpha > 0$, will the basic gradient descent algorithm terminate at a minimum?

$$f'(\theta) = 0$$

- A. Never B. Maybe C. Yes with enough iterations

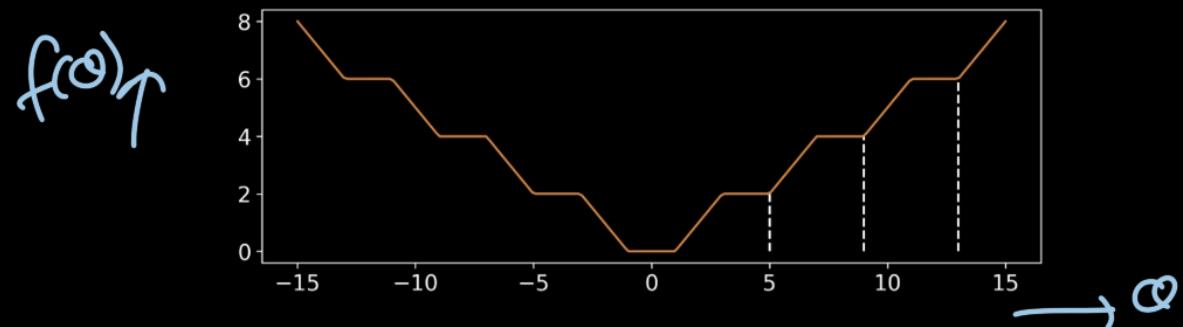
- (b) Assuming that θ starts in a sloped region and $\alpha > 0$, will the basic gradient descent algorithm find the minimum?

- A. Never B. Maybe C. Yes with enough iterations



Question:

13. Consider the following function of $f(\theta)$, which alternates between completely flat regions and regions of absolute slope equal to 1. **There is only one correct answer for each part.**



- (a) Assuming that θ starts in a flat region that is not a minimum and $\alpha > 0$, will the basic gradient descent algorithm terminate at a minimum?

$$f'(\theta) = 0$$

- A. Never B. Maybe C. Yes with enough iterations

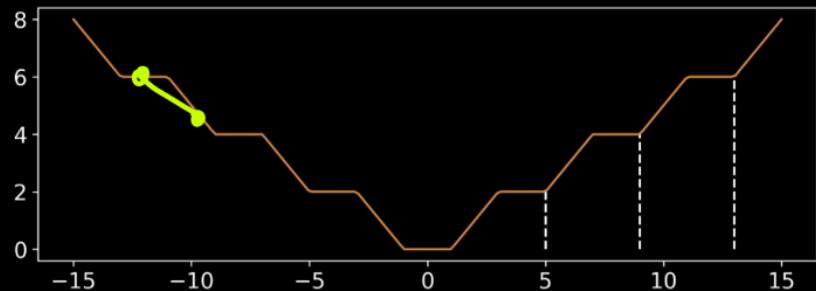
- (b) Assuming that θ starts in a sloped region and $\alpha > 0$, will the basic gradient descent algorithm find the minimum?

- A. Never B. Maybe C. Yes with enough iterations



Question:

13. Consider the following function of $f(\theta)$, which alternates between completely flat regions and regions of absolute slope equal to 1. **There is only one correct answer for each part.**



Is $f(\theta)$ convex?

- A. Yes
- B. No ✓

CLASSES



(a) Assuming that θ starts in a flat region that is not a minimum and $\alpha > 0$, will the basic gradient descent algorithm terminate at a minimum? Note that the basic gradient descent algorithm is just the same as version with momentum on the previous page, but where $\gamma = 0$.

- A. Never B. Maybe C. Yes with enough iterations

(b) Assuming that θ starts in a sloped region and $\alpha > 0$, will the basic gradient descent algorithm find the minimum?

- A. Never B. Maybe C. Yes with enough iterations

Is $f(\theta)$ convex?

- A. Yes
 B. No



Question:

27. [10 Pts] Consider the following loss function based on data x_1, \dots, x_n with mean \bar{x} :

$$\ell(\beta) = \log \beta + \frac{\bar{x}}{\beta} + \frac{1}{n} \sum_{i=1}^n e^{-x_i/\beta}$$

Given an estimate $\beta^{(t)}$, write out the update $\beta^{(t+1)}$ after one iteration of gradient descent with step size α . Show your work in the box below.



Solution: The update is

$$\beta^{(t+1)} \leftarrow \beta^{(t)} - \alpha \ell'(\beta^{(t)}),$$

where

$$\begin{aligned}\ell'(\beta) &= \frac{1}{\beta} \left(1 - \frac{\bar{x}}{\beta} + \frac{1}{n\beta} \sum_{i=1}^n x_i e^{-x_i/\beta} \right) \\ &= \frac{1}{\beta} - \frac{\bar{x}}{\beta^2} + \frac{1}{n\beta^2} \sum_{i=1}^n x_i e^{-x_i/\beta}\end{aligned}$$

Alternate notation:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} - \alpha \frac{\partial \ell}{\partial \beta} \Big|_{\beta=\beta^{(t)}}$$

With everything substituted in:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} - \alpha \left(\frac{1}{\beta^{(t)}} - \frac{\bar{x}}{\beta^{(t)2}} + \frac{1}{n\beta^{(t)2}} \sum_{i=1}^n x_i e^{-x_i/\beta^{(t)}} \right)$$



Question:

(m) [3 pts] Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous, smooth function whose derivative $f'(x)$ is also continuous. Suppose f has a unique global minimum $x^* \in (-\infty, \infty)$, and you are using **gradient descent** to find x^* . You fix some $x^{(0)} \in \mathbb{R}$ and $\epsilon > 0$, and run $x^{(t)} = x^{(t-1)} - \epsilon f'(x^{(t-1)})$ repeatedly. Which of the following statements are true?

- Gradient descent is sure to converge, to *some* value, for any step size $\epsilon > 0$
 - If f has a local minimum x' different from the global one, i.e., $x' \neq x^*$, and $x^{(t)} = x'$ for some t , gradient descent will not converge to x^*
 - Assuming gradient descent converges, it converges to x^* if and only if f is convex
- {
- If, additionally, f is the objective function of logistic regression, and gradient descent converges, then it converges to x^*

skip this option



(m) [3 pts] Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous, smooth function whose derivative $f'(x)$ is also continuous. Suppose f has a unique global minimum $x^* \in (-\infty, \infty)$, and you are using **gradient descent** to find x^* . You fix some $x^{(0)} \in \mathbb{R}$ and $\epsilon > 0$, and run $x^{(t)} = x^{(t-1)} - \epsilon f'(x^{(t-1)})$ repeatedly. Which of the following statements are true?

Gradient descent is sure to converge, to *some* value, for any step size $\epsilon > 0$

If f has a local minimum x' different from the global one, i.e., $x' \neq x^*$, and $x^{(t)} = x'$ for some t , gradient descent will not converge to x^*

Assuming gradient descent converges, it converges to x^* if and only if f is convex

If, additionally, f is the objective function of logistic regression, and gradient descent converges, then it converges to x^*

The top-left option is false because for a large enough step size, gradient descent may not converge. The bottom-left option is correct because $f'(x') = 0$, so gradient descent will never move from a local minimum. The top-right option is false because you could “accidentally” initialize GD at x^* even if f is non-convex. The bottom-right option is correct because the objective of logistic regression is convex.



Question:

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous, smooth function whose derivative $f'(x)$ is also continuous. Suppose f has a unique global minimum x^* , and you are using gradient descent to find x^* . You fix some x_0 and $\eta > 0$, and apply $x_{t+1} = x_t - \eta f'(x_t)$ repeatedly. Which of the following statements is true?

- (A) Gradient descent is sure to converge to some value, for any step size $\eta > 0$.
- (B) Assuming gradient descent converges, it converges to x^* if and only if f is convex.
- (C) If f has a local minimum x' different from the global one, and $x_t = x'$ for some t , gradient descent will not converge to x^* .



Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous, smooth function whose derivative $f'(x)$ is also continuous. Suppose f has a unique global minimum x^* , and you are using gradient descent to find x^* . You fix some x_0 and $\eta > 0$, and apply $x_{t+1} = x_t - \eta f'(x_t)$ repeatedly. Which of the following statements is true?

- (A) Gradient descent is sure to converge to some value, for any step size $\eta > 0$.
- (B) Assuming gradient descent converges, it converges to x^* if and only if f is convex.
- (C) If f has a local minimum x' different from the global one, and $x_t = x'$ for some t , gradient descent will not converge to x^* .

Solution:

The solution is (C). For a large enough step size, gradient descent may not converge. If $x_0 = x^*$ then GD converges regardless of convexity. When $x_t = x'$, then $f'(x_t) = 0$, so it converges to x' instead of x^* .



Question:

Gradient descent is sure to converge, to some value, for any step size greater than 0.

Solution: False. The step-size can be too large, which may cause the algorithm to diverge.





Question:

Suppose that we are given $f(x) = x^3 + x^2$ and learning rate $\alpha = 1/4$.

- a) [5 Points] First of all, write down the updating rule for gradient descent in general and for this function.
- b) [5 Points] If we start at $x_0 = -1$, should we go left or right? Can you verify mathematically? What is x_1 ? Can gradient descent converge? If so, where it might converge to, given appropriate step size?
- c) [5 Points] If we start at $x_0 = 1$, should we go left or right? Can you verify mathematically? What is x_1 ? Can gradient descent converge? If so, where it might converge to, given appropriate step size?
- d) [5 Points] Write down 1 condition to terminate the gradient descent algorithm (in general).



a) [5 Points] First of all, write down the updating rule for gradient descent in general and for this function.

Solution: In general, the updating rule for gradient descent is:

$$x_{i+1} = x_i - \alpha \nabla f(x_i) = x_i - \alpha \frac{\partial f}{\partial x}(x_i),$$

where $\alpha \in \mathbb{R}_+$ is the learning rate or step size. For this function, since f is a single-variable function, we can write down the updating rule as:

$$x_{i+1} = x_i - \alpha \frac{df}{dx}(x_i) = x_i - \alpha f'(x_i).$$

We also have:

$$\frac{df}{dx} = f'(x) = 3x^2 + 2x,$$

thus the updating rule can be written down as:

$$x_{i+1} = x_i - \alpha(3x_i^2 + 2x_i) = -\frac{3}{4}x_i^2 + \frac{1}{2}x_i.$$

- b) [5 Points] If we start at $x_0 = -1$, should we go left or right? Can you verify mathematically? What is x_1 ? Can gradient descent converge? If so, where it might converge to, given appropriate step size?

Solution: We have

$$f'(x_0) = f'(-1) = 3(-1)^2 + 2(-1) = 1 > 0,$$

so we go left, and

$$x_1 = x_0 - \alpha f'(x_0) = -1 - \frac{1}{4} = -\frac{5}{4}.$$

Intuitively, the gradient descent cannot converge in this case because

$$\lim_{x \rightarrow -\infty} f(x) = -\infty,$$

We need to find all local minimums and local maximums. First, we solve the equation $f'(x) = 0$ to find all critical points. We have:

$$f'(x) = 0 \Leftrightarrow 3x^2 + 2x = 0 \Leftrightarrow x = -\frac{2}{3} \text{ and } x = 0.$$

Now, we consider the second-order derivative:

$$f''(x) = \frac{d^2 f}{dx^2} = 6x + 2.$$

We have $f''(x) = 0$ only when $x = -1/3$. Thus, for $x < -1/3$, $f''(x)$ is negative or the slope $f'(x)$ decreases; and for $x > -1/3$, $f''(x)$ is positive or the slope $f'(x)$ increases. Keep in mind that $-1 < -2/3 < -1/3 < 0 < 1$. Therefore, f has a local maximum at $x = -2/3$ and a local minimum at $x = 0$. If the gradient descent starts at $x_0 = -1$ and it always goes left then it will never meet the local minimum at $x = 0$, and it will go left infinitely. We say the gradient descent cannot converge, or is divergent.



- c) [5 Points] If we start at $x_0 = 1$, should we go left or right? Can you verify mathematically? What is x_1 ? Can gradient descent converge? If so, where it might converge to, given appropriate step size?

Solution: We have

$$f'(x_0) = f'(-1) = 3 \cdot 1^2 + 2 \cdot 1 = 5 > 0,$$

so we go left, and

$$x_1 = x_0 - \alpha f'(x_0) = 1 - \frac{1}{4} \cdot 5 = -\frac{1}{4}.$$

From the previous part, function f has a local minimum at $x = 0$, so the gradient descent can converge (given appropriate step size) at this local minimum.

- d) [5 Points] Write down 1 condition to terminate the gradient descent algorithm (in general).

Solution: There are several ways to terminate the gradient descent algorithm:

- If the change in the optimization objective is too small, i.e. $|f(x_i) - f(x_{i+1})| < \epsilon$ where ϵ is a small constant,
- If the gradient is close to zero or the norm of the gradient is very small, i.e. $\|\nabla f(x_i)\| < \lambda$ where λ is a small constant.



MSQ Question:

Which of the following are true about gradient descent? (select all statements that are true.)

- After each iteration, we modify the weight vector in the direction of the gradient.
- We have to choose a non-variable learning rate.
- After each iteration, we modify the weight vector in the direction of the negative gradient.
- In the gradient descent algorithm each update of the weight vector depends on all the training examples.

Question 3

1 / 1 pts

Which of the following are true about gradient descent? (select all statements that are true.)

After each iteration, we modify the weight vector in the direction of the gradient.

We have to choose a non-variable learning rate.

After each iteration, we modify the weight vector in the direction of the negative gradient.

In the gradient descent algorithm each update of the weight vector depends on all the training examples.

SES

Correct!

Correct!

**Question:**

Select all that apply: Consider the convex function $f(z) = z^2$. Let α be our learning rate in gradient descent.

For which values of α will $\lim_{t \rightarrow \infty} f(z^{(t)}) = 0$, assuming the initial value of z is $z^{(0)} = 1$ and $z^{(t)}$ is the value of z after the t -th iteration of gradient descent?

- $\alpha = 0$
- $\alpha = \frac{1}{2}$
- $\alpha = 1$
- $\alpha = 2$
- None of the above

Numerical answer: Give the range of all values for $\alpha \geq 0$ such that $\lim_{t \rightarrow \infty} f(z^{(t)}) = 0$, assuming the initial value of z is $z^{(0)} = 1$.



Select all that apply: Consider the convex function $f(z) = z^2$. Let α be our learning rate in gradient descent.

For which values of α will $\lim_{t \rightarrow \infty} f(z^{(t)}) = 0$, assuming the initial value of z is $z^{(0)} = 1$ and $z^{(t)}$ is the value of z after the t -th iteration of gradient descent?

- $\alpha = 0$
- $\alpha = \frac{1}{2}$ $\alpha = \frac{1}{2}$
- $\alpha = 1$
- $\alpha = 2$
- None of the above

Numerical answer: Give the range of all values for $\alpha \geq 0$ such that $\lim_{t \rightarrow \infty} f(z^{(t)}) = 0$, assuming the initial value of z is $z^{(0)} = 1$.

(0, 1).





Question:

Convexity is a desirable property in machine learning because it:

- (a) guarantees gradient descent finds a global minimum in optimization problems for functions that have a global minimum
- (b) helps to avoid the model overfitting
- (c) speeds up model training
- (d) reduces model complexity

Skip these options



Convexity is a desirable property in machine learning because it:

- (a) guarantees gradient descent finds a global minimum in optimization problems for functions that have a global minimum
- (b) helps to avoid the model overfitting
- (c) speeds up model training
- (d) reduces model complexity

Correct answers: (a)

✓

$$x_{\text{new}} = x_{\text{old}} - \alpha \nabla f(x)$$

 GO
CLASSES



Constrained optimisation

