



AI VIET NAM

@aivietnam.edu.vn

A Bridge to Linear Regression

Quang-Vinh Dinh
PhD in Computer Science

Objective

Feature		Label	
	area	price	
	6.7	9.1	
	4.6	5.9	
	3.5	4.6	
	5.5	6.7	

House price data

Features			Label
TV	↕ Radio	↕ Newspaper	↕ Sales
230.1	37.8	69.2	22.1
44.5	39.3	45.1	10.4
17.2	45.9	69.3	12
151.5	41.3	58.5	16.5
180.8	10.8	58.4	17.9

Advertising data

if area=6.0, price=?

if TV=55.0, Radio=34.0,
and Newspaper=62.0,
price=?

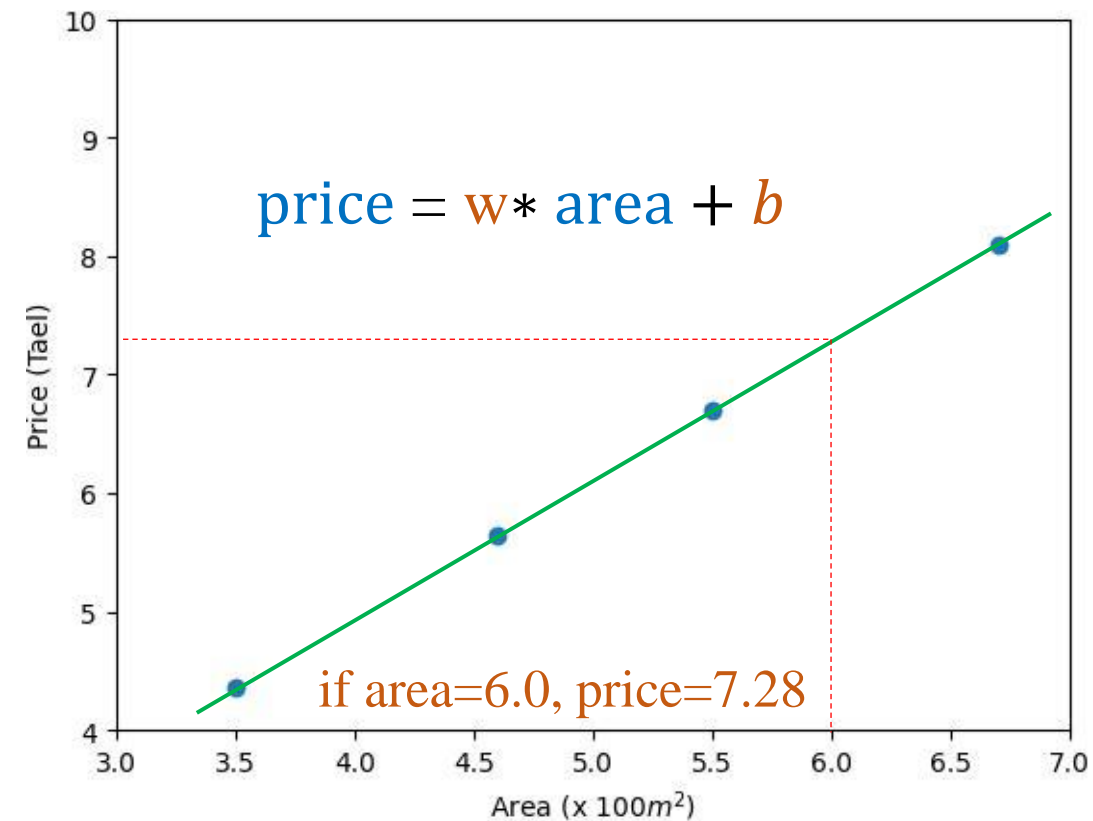
Features														Label
crim	↕ zn	↕ indus	↕ chas	↕ nox	↕ rm	↕ age	↕ dis	↕ rad	↕ tax	↕ ptratio	↕ black	↕ lstat	↕ medv	
0.00632	18	2.31	0	0.538	6.575	65.2	4.09	1	296	15.3	396.9	4.98	24	
0.02731	0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.9	9.14	21.6	
0.03237	0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4	
0.06905	0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.9	5.33	36.2	
0.08829	12.5	7.87	0	0.524	6.012	66.6	5.5605	5	311	15.2	395.6	12.43	22.9	

Boston House Price Data

House Price Prediction

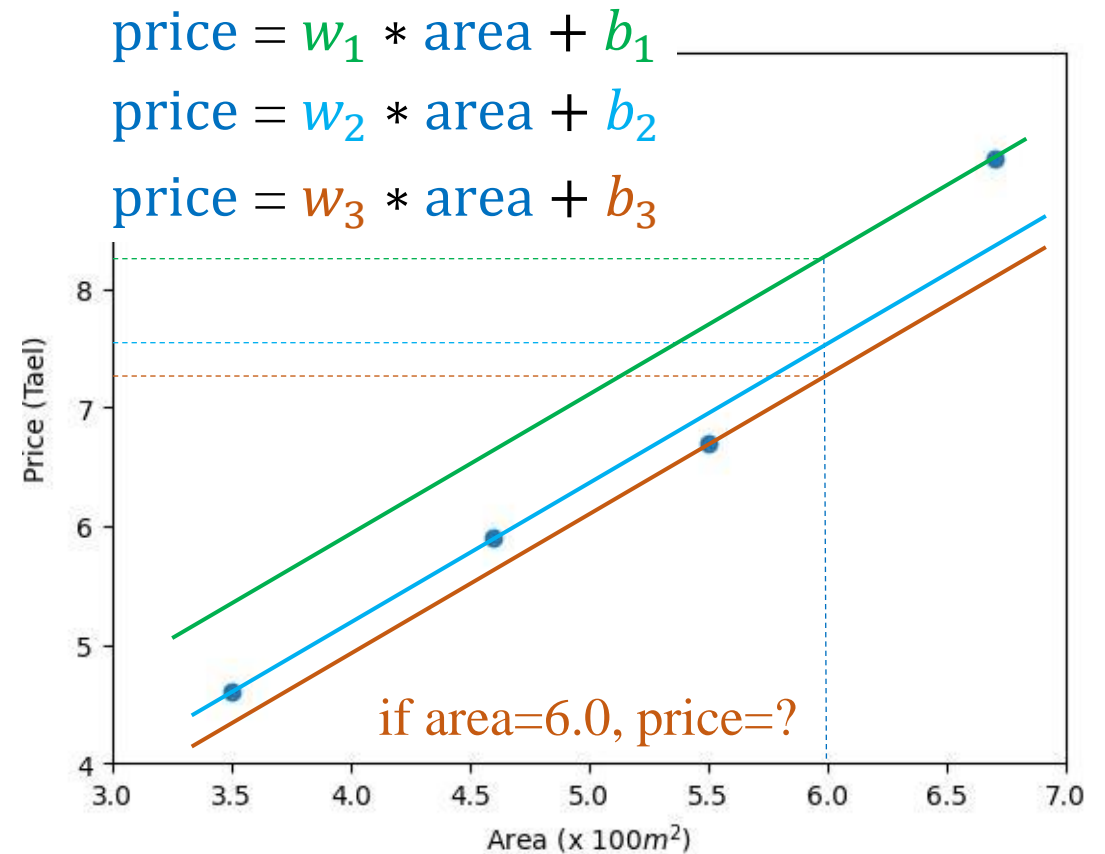
House price data

Feature	Label
area	price
6.7	8.1
4.6	5.6
3.5	4.3
5.5	6.7



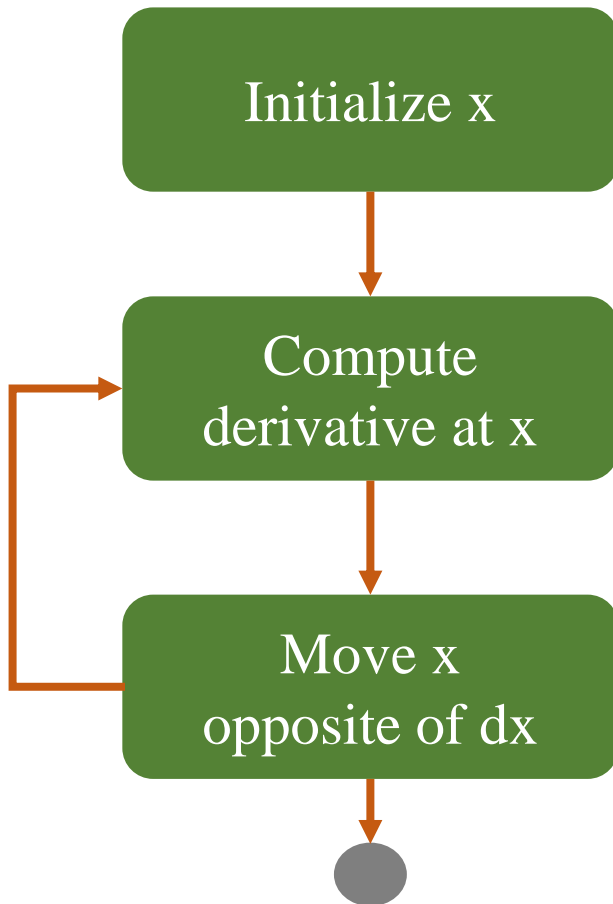
Feature	Label
area	price
6.7	9.1
4.6	5.9
3.5	4.6
5.5	6.7

House price data

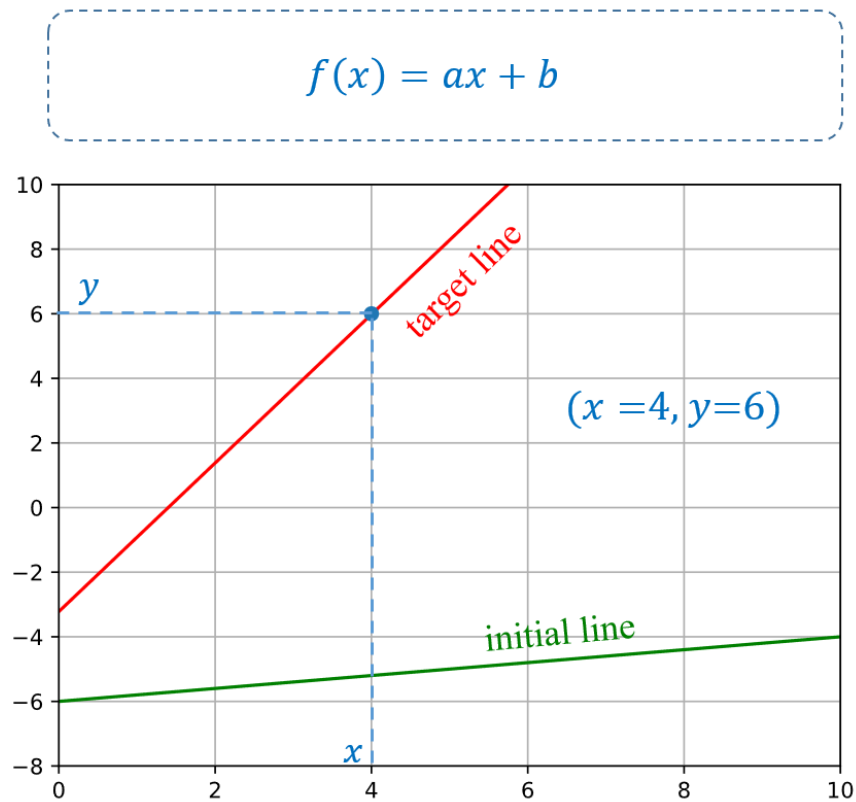


Objectives

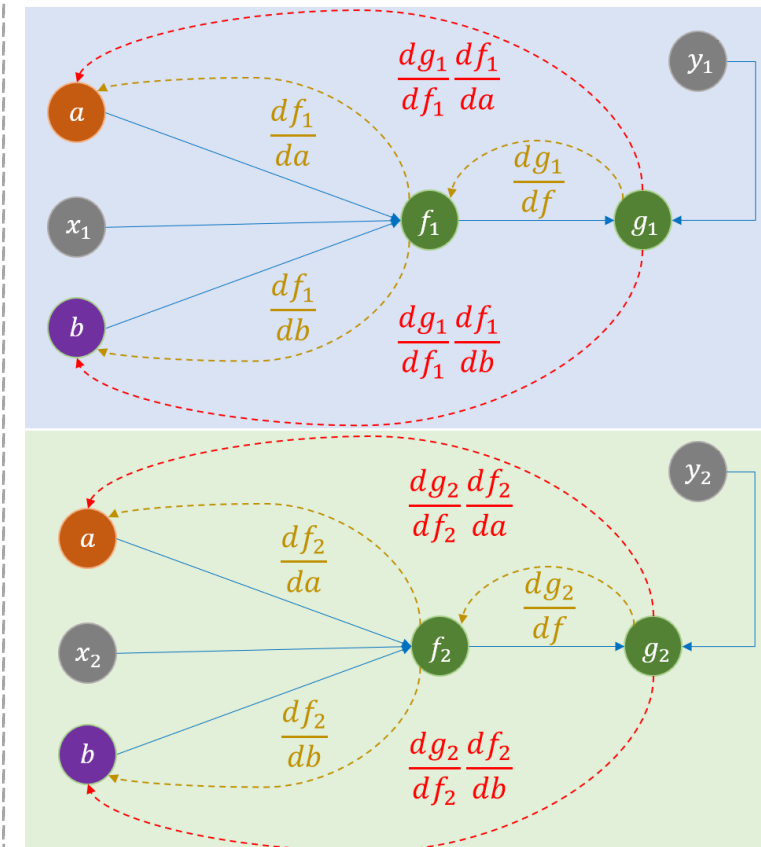
Optimization



Problem Solving



Toward Linear Regression



Outline

SECTION 1

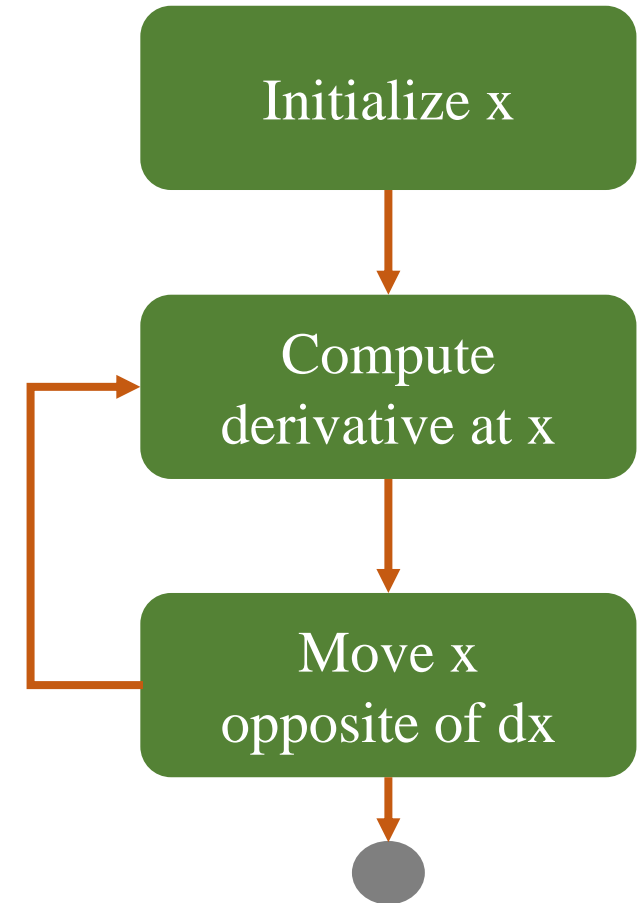
Optimization

SECTION 2

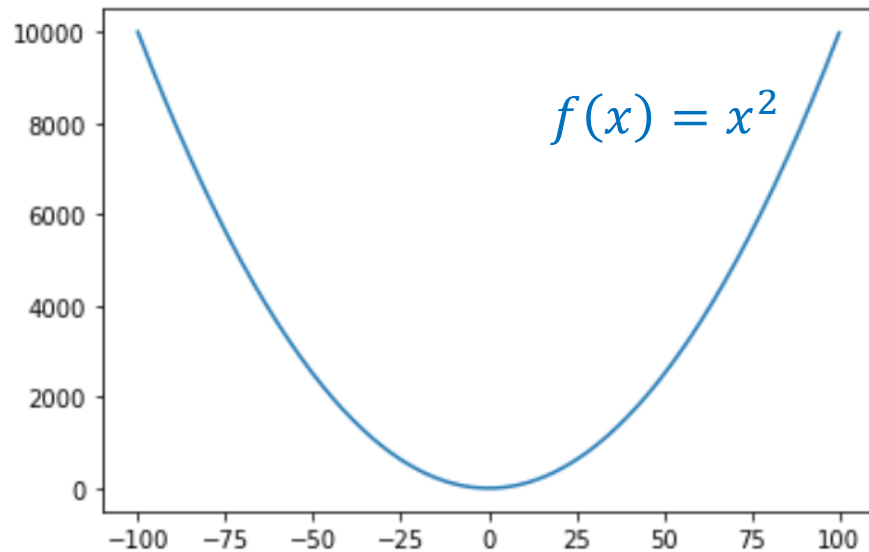
Problem Solving

SECTION 3

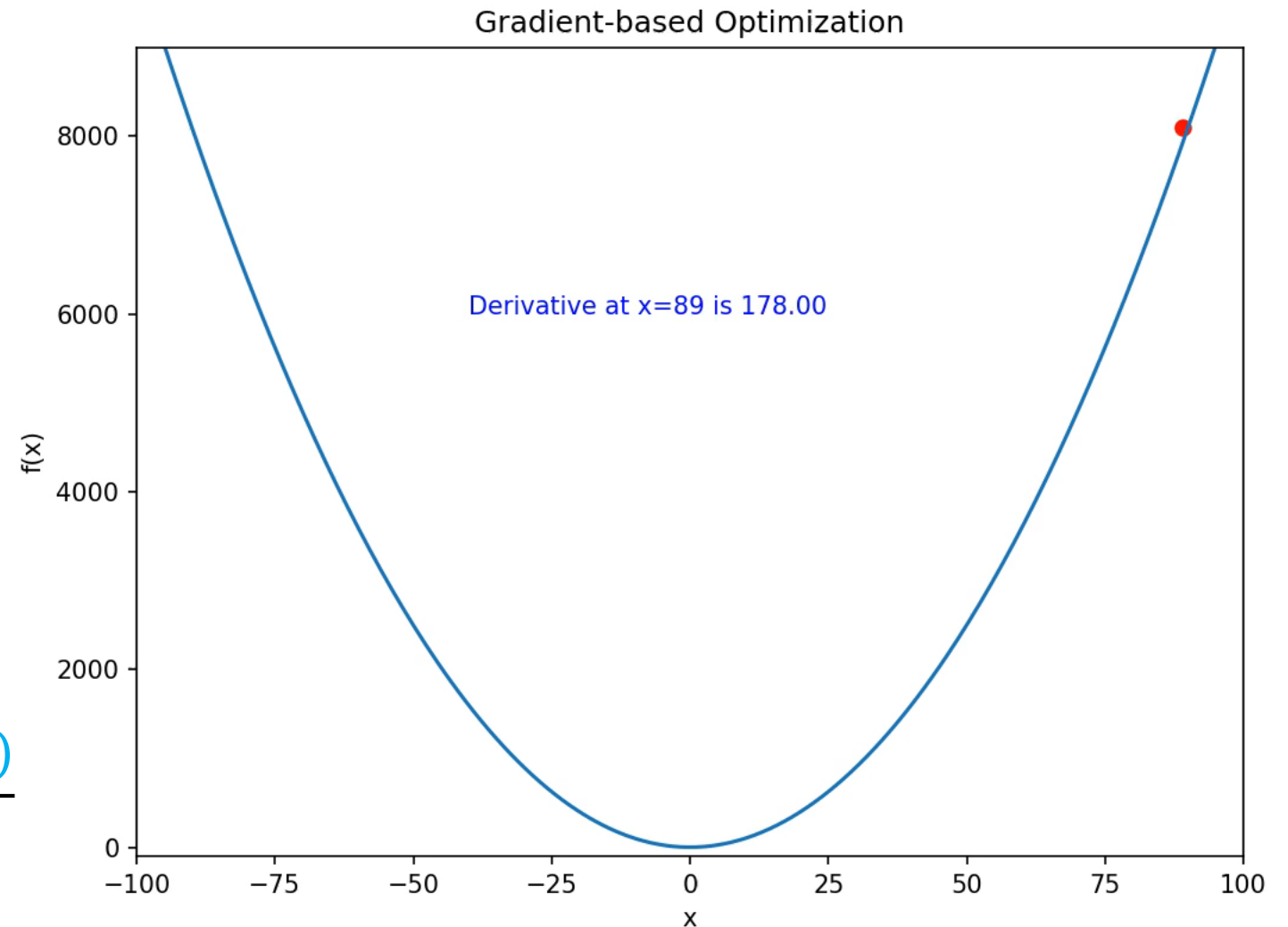
Towards Linear Regression



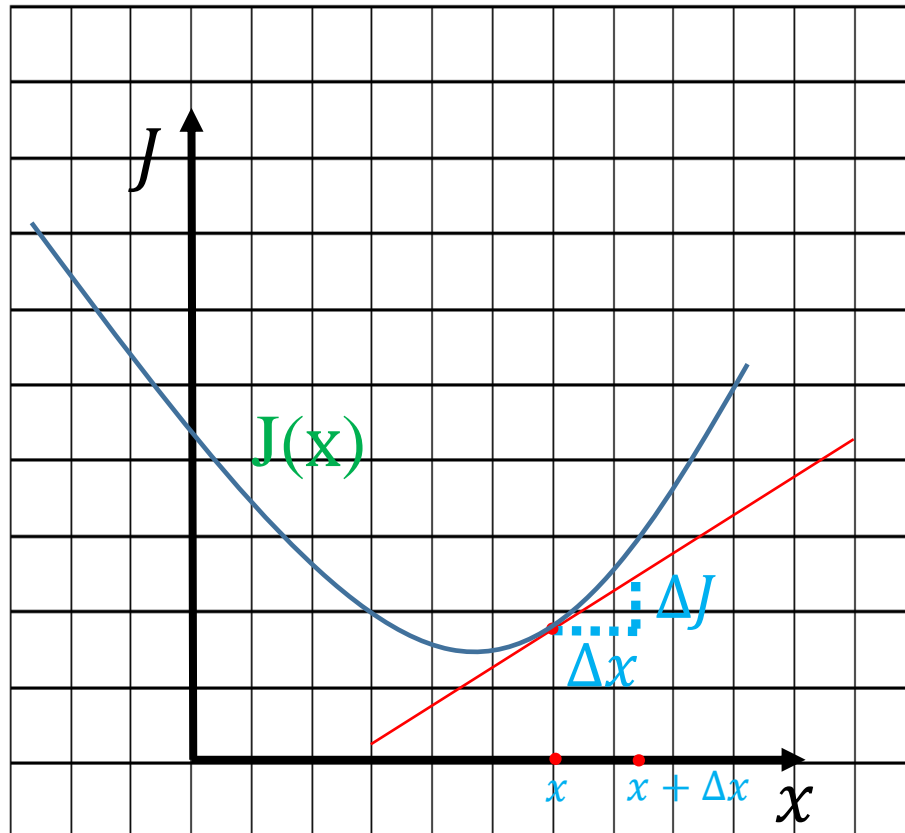
❖ Square function



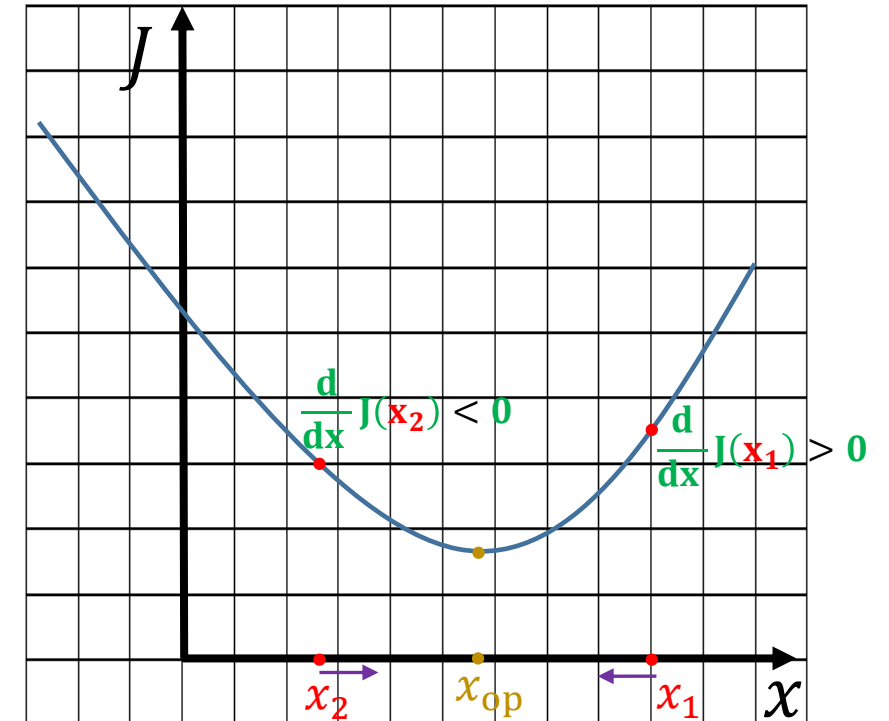
$$\frac{d}{dx} f(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$



❖ Gradient descent



$$\frac{d}{dx}J(x) = \lim_{\Delta x \rightarrow 0} \frac{J(x + \Delta x) - J(x)}{\Delta x}$$



$$x_{new} = x_{old} - \eta \frac{d}{dx}J(x_{old})$$

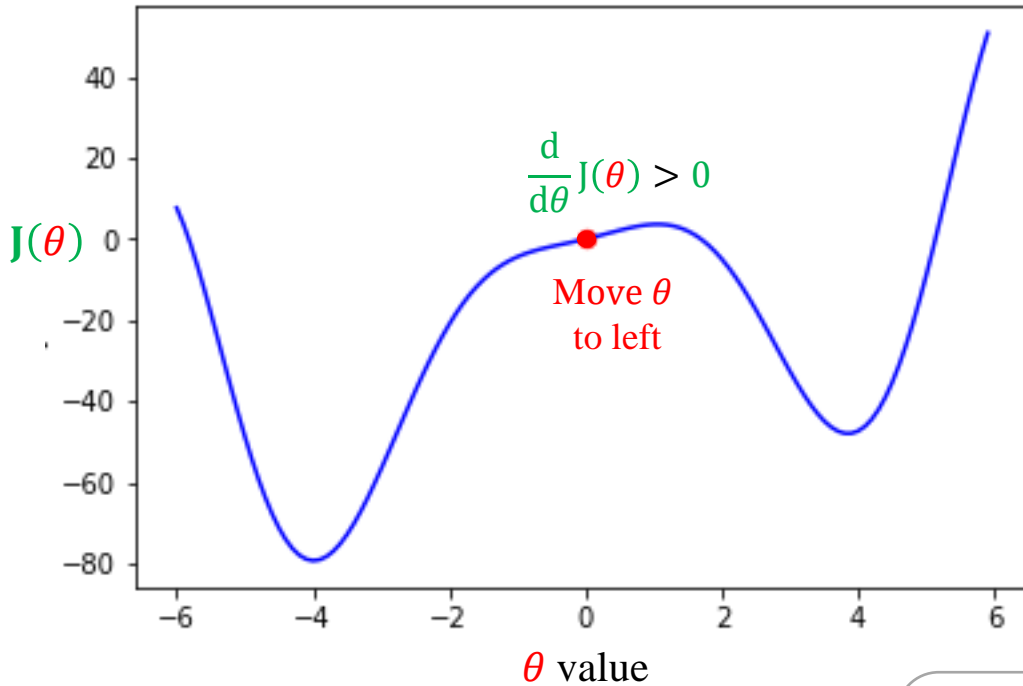
Derivate at x_{old}

learning rate

Optimization

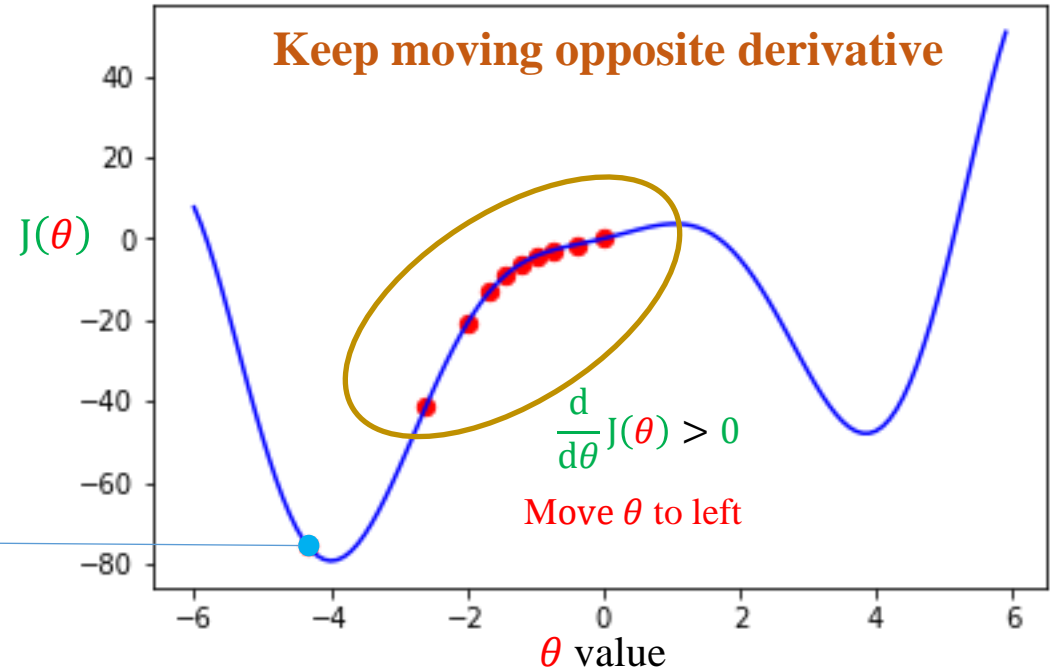
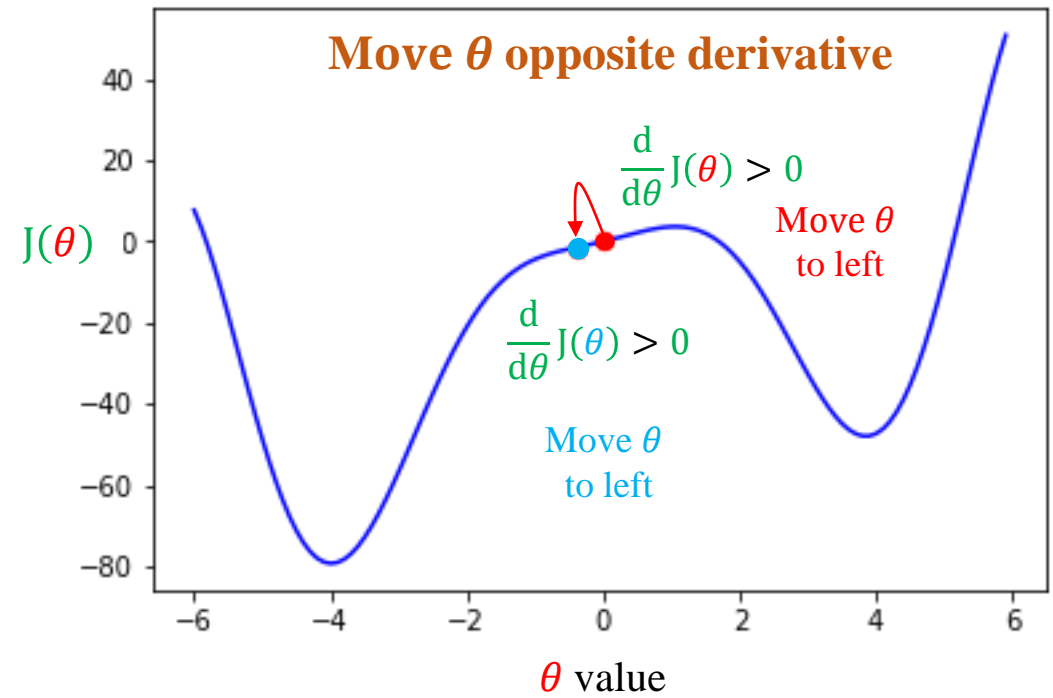
❖ Gradient descent

Initialize θ



$\frac{d}{d\theta} J(\theta) < 0$
Move θ to right

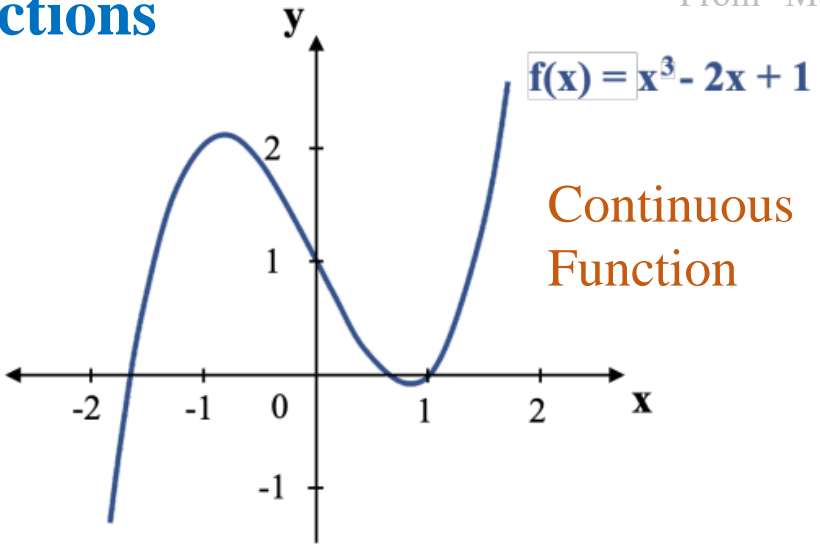
Move θ opposite derivative



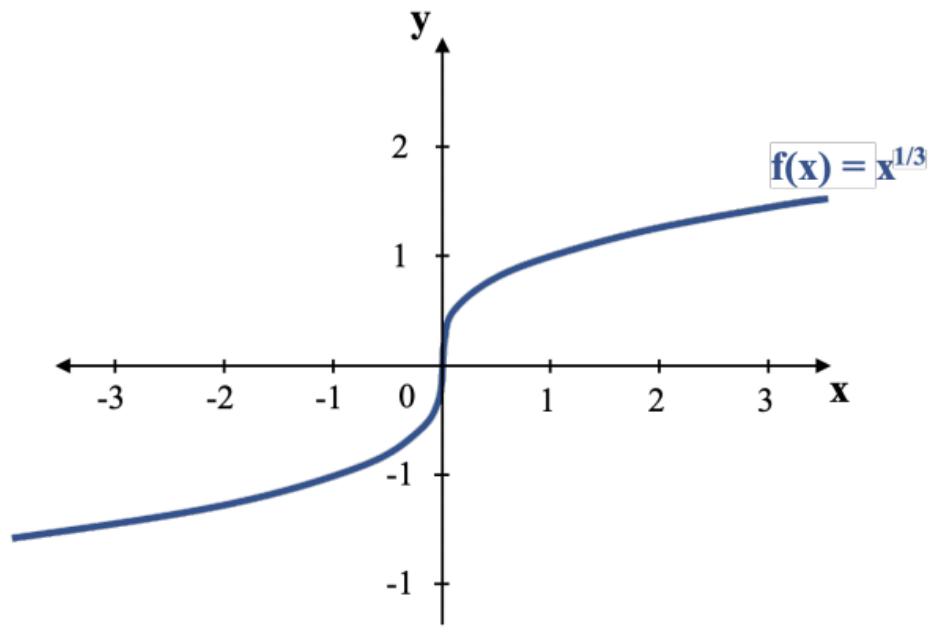
Optimization Algorithms

Loss functions

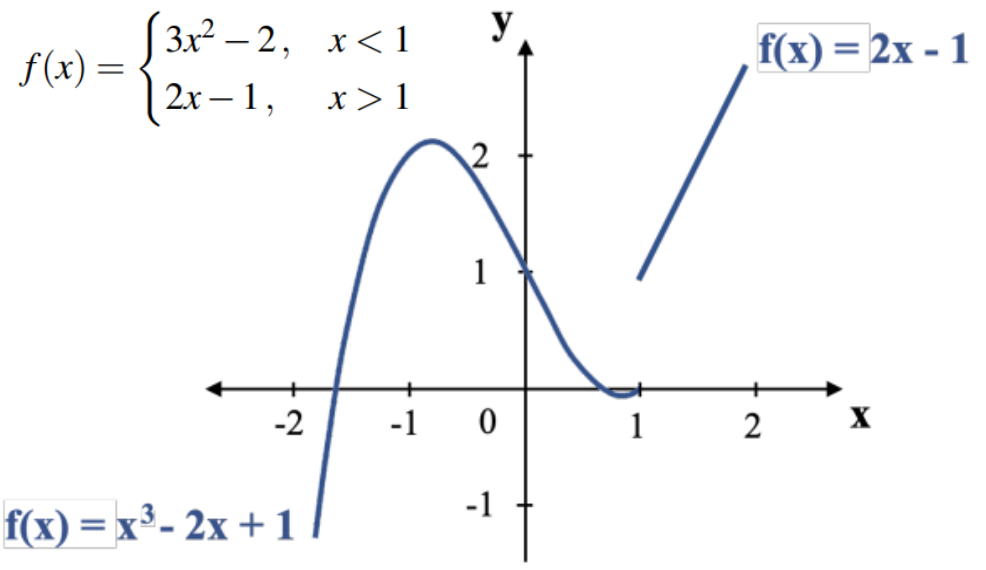
From “Machine Learning Simplified”



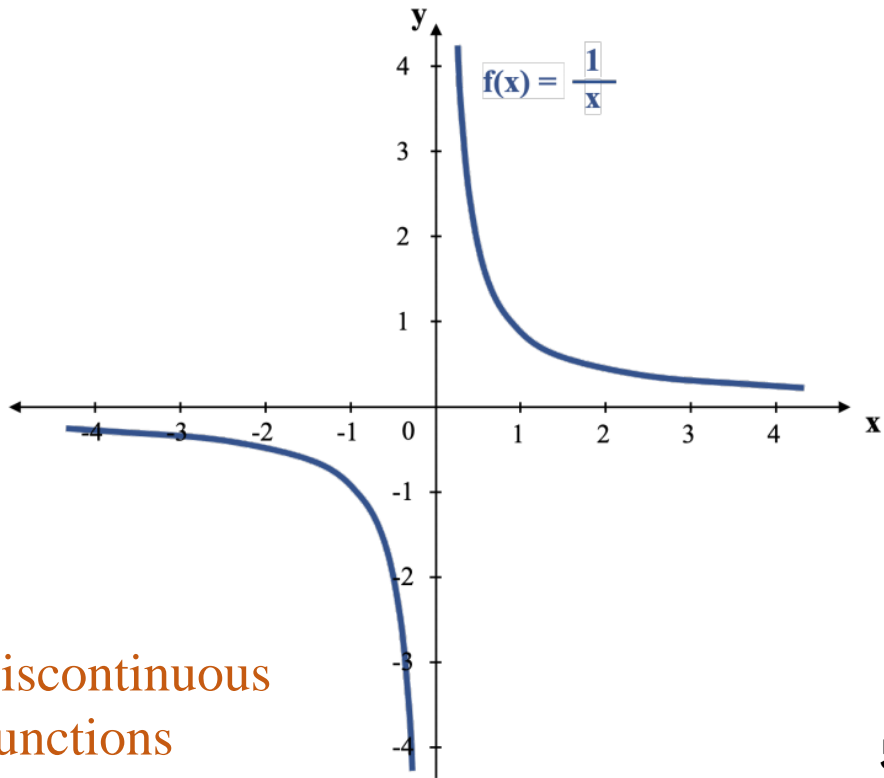
Continuous
Function



Continuous non-differentiable functions

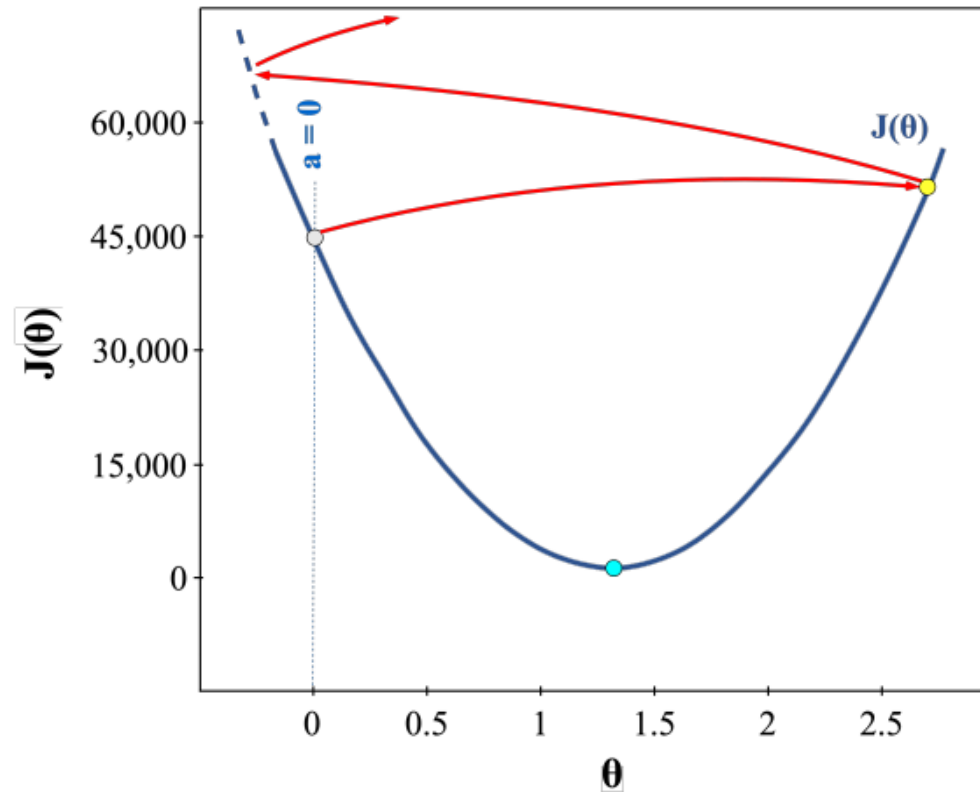


$f(x) = x^3 - 2x + 1$

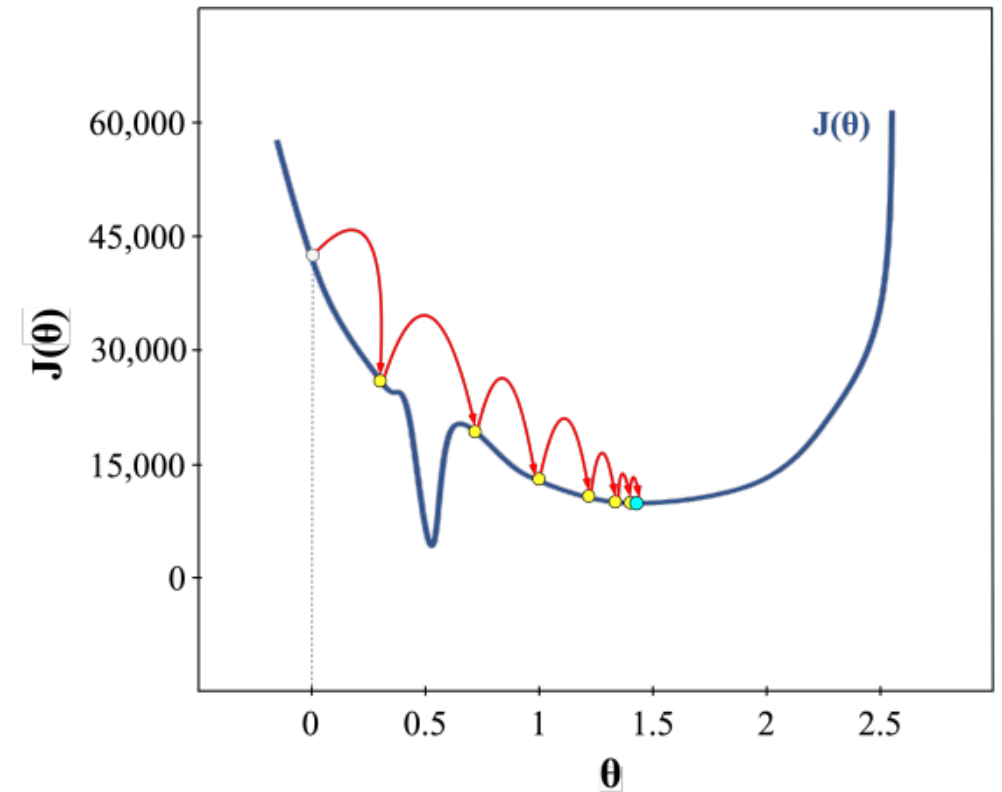


Discontinuous
Functions

❖ Learning rate



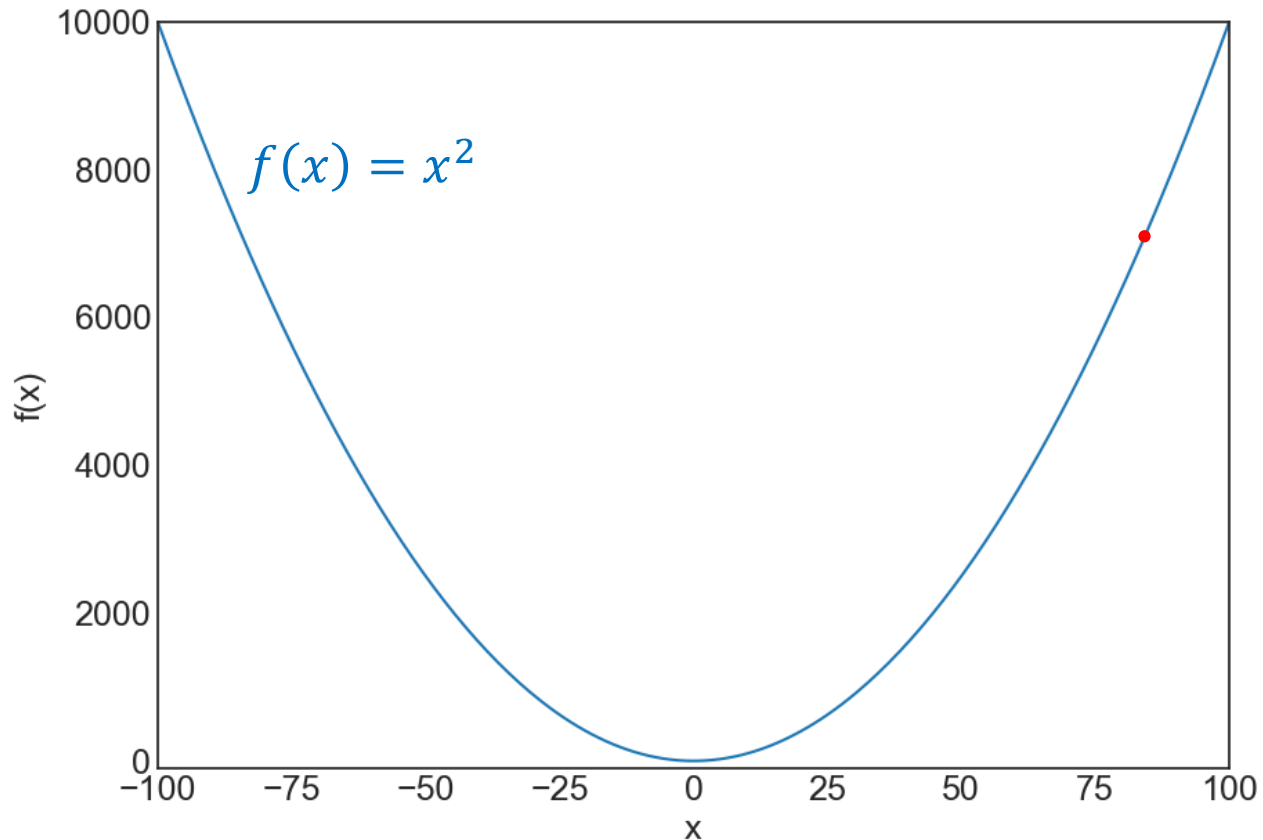
(a) Gradient descent missing global minimum on a convex cost function due to a very large learning rate.



(b) Gradient Descent missing global minimum on a non-convex cost function due to a very large learning rate.

Observation 1

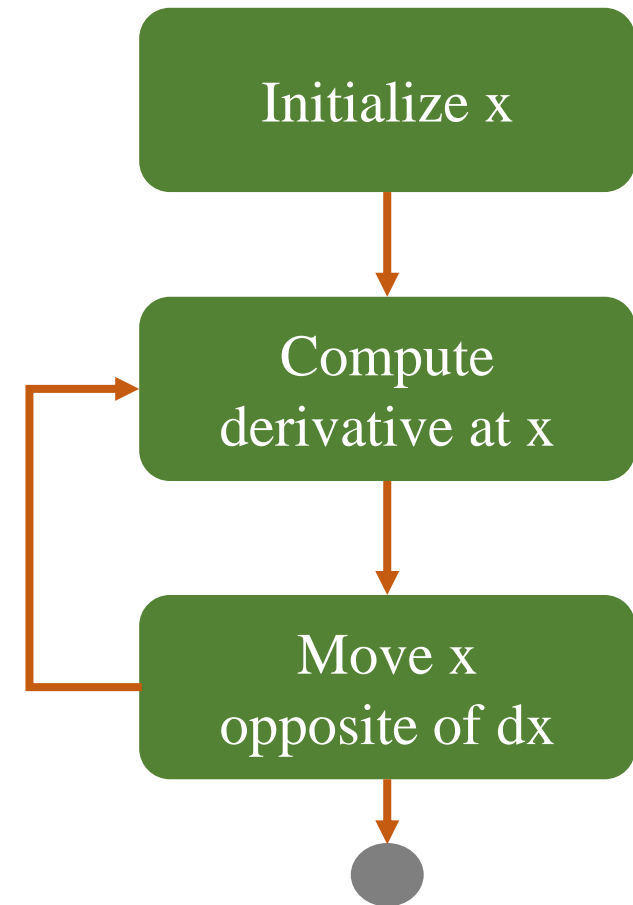
❖ Square function



$$-100 \leq x \leq 100$$

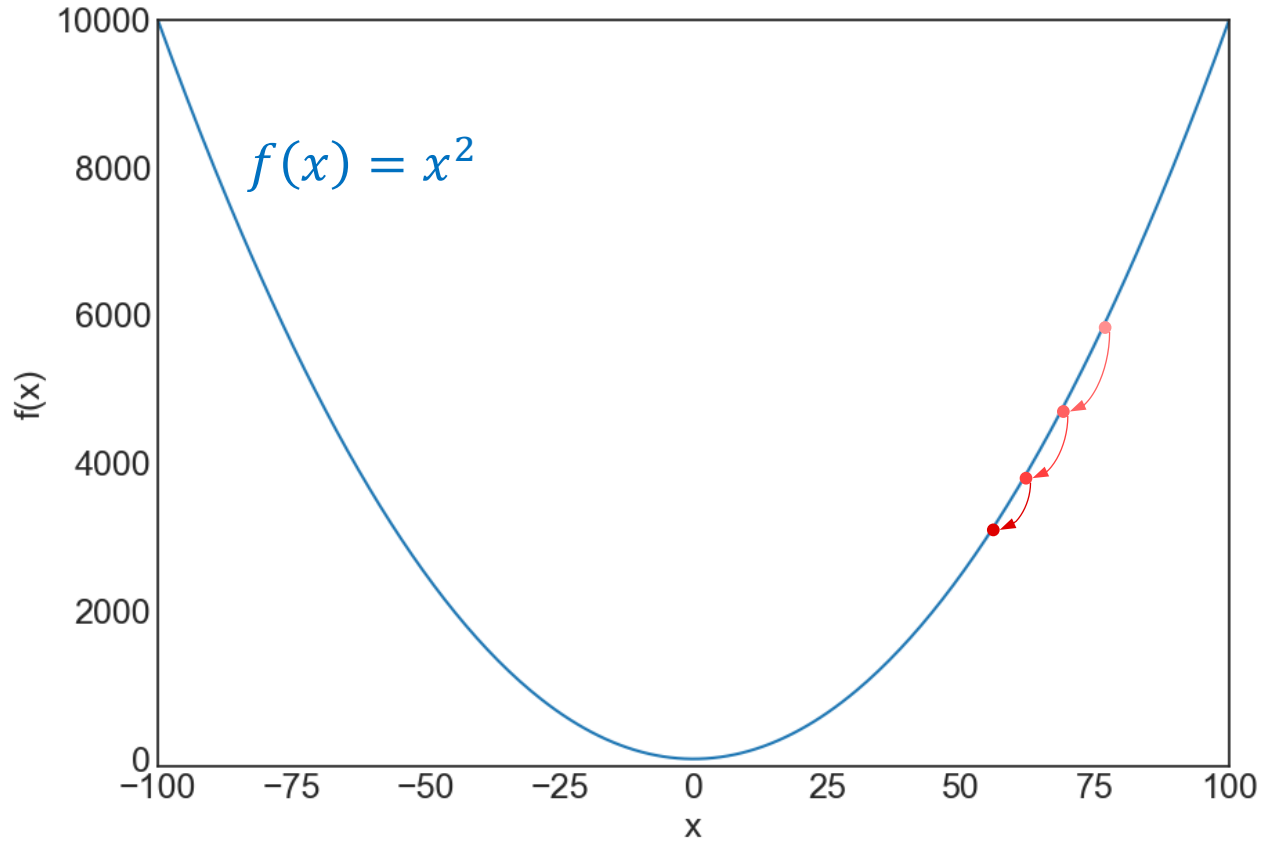
$$x \in \mathbb{N}$$

$$x_t = x_{t-1} - \eta f'(x_{t-1})$$



Optimization

❖ Square function



$$\begin{aligned} -100 \leq x \leq 100 \\ x \in \mathbb{N} \end{aligned}$$

$$x_t = x_{t-1} - \eta f'(x_{t-1})$$

$$x_0 = 70.0 \quad \eta = 0.1$$

$$f'(x_0) = 140.0$$

$$x_1 = x_0 - \eta f'(x_0) = 56.0$$

$$f'(x_1) = 112.0$$

$$x_2 = x_1 - \eta f'(x_1) = 44.8$$

$$f'(x_2) = 89.6$$

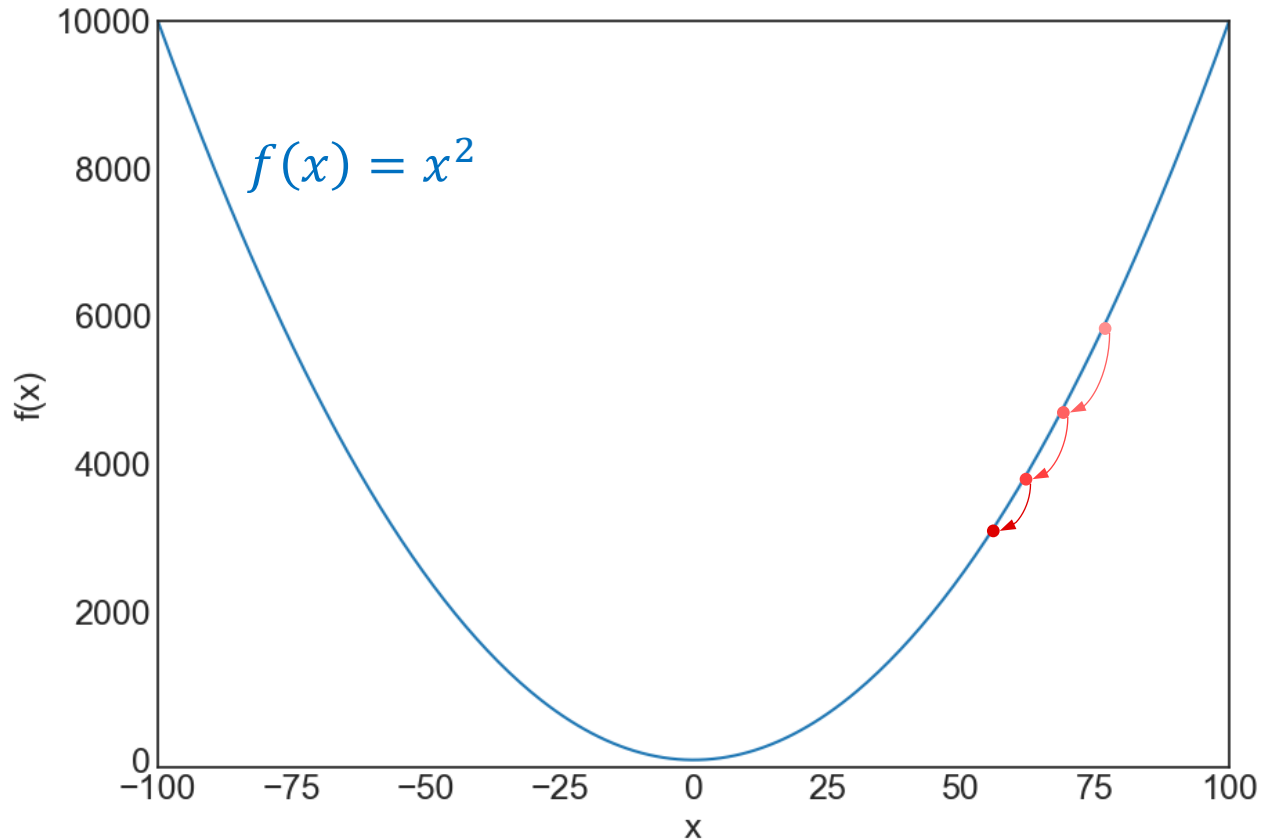
$$x_3 = x_2 - \eta f'(x_2) = 35.84$$

$$f'(x_3) = 71.68$$

$$x_4 = x_3 - \eta f'(x_3) = 28.672$$

Optimization

❖ Square function



Keep doing

$$x_t = x_{t-1} - \eta f'(x_{t-1})$$

$$x_{10} = 6.012 \quad \eta = 0.1$$

$$f'(x_{10}) = 12.02$$

$$x_{11} = x_{10} - \eta f'(x_{10}) = 4.81$$

$$f'(x_{11}) = 9.62$$

$$x_{12} = x_{11} - \eta f'(x_{11}) = 3.84$$

$$f'(x_{12}) = 7.69$$

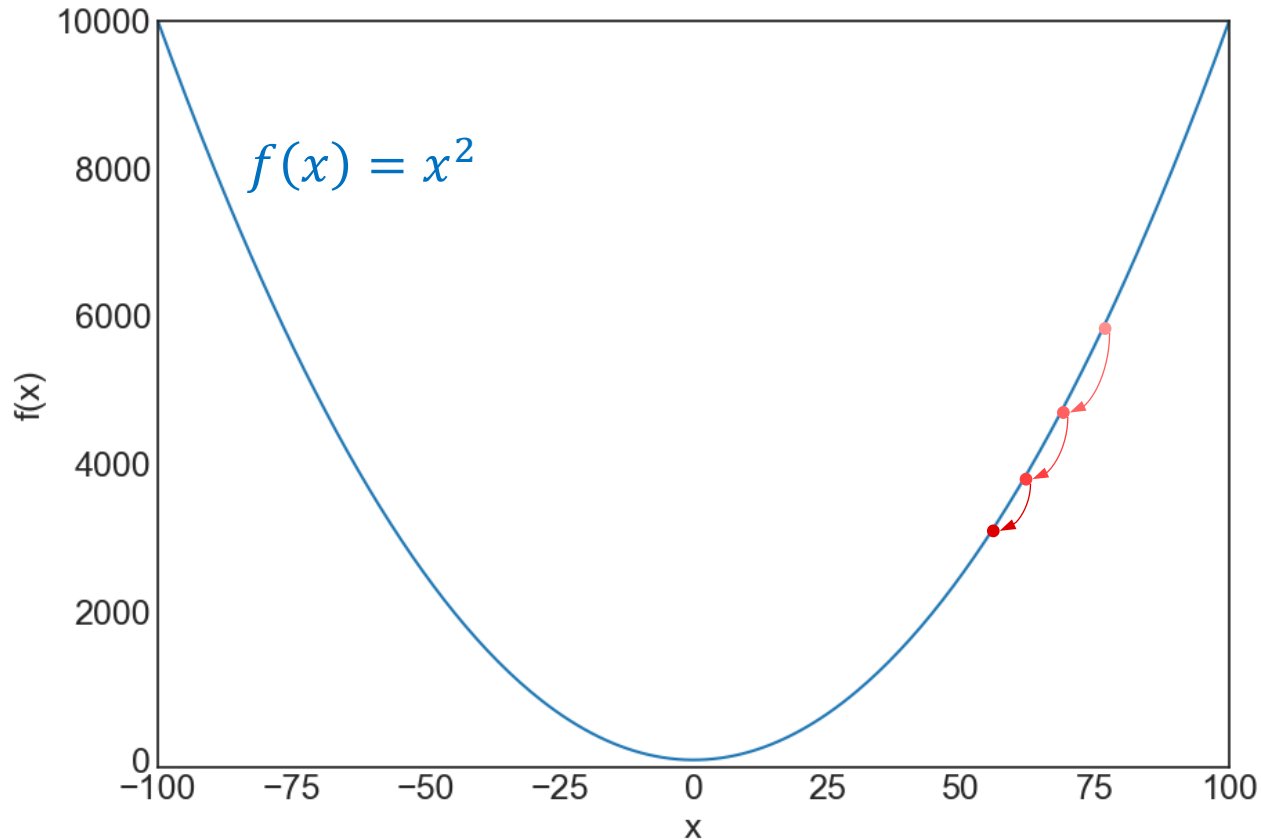
$$x_{13} = x_{12} - \eta f'(x_{12}) = 3.078$$

$$f'(x_{13}) = 6.15$$

$$x_{14} = x_{13} - \eta f'(x_{13}) = 2.46$$

Optimization

❖ Square function



Keep doing

$$x_t = x_{t-1} - \eta f'(x_{t-1})$$

$$x_{30} = 0.069 \quad \eta = 0.1$$

$$f'(x_{30}) = 0.138$$

$$x_{31} = x_{30} - \eta f'(x_{30}) = 0.055$$

$$f'(x_{31}) = 0.11$$

$$x_{32} = x_{31} - \eta f'(x_{31}) = 0.044$$

$$f'(x_{32}) = 0.88$$

$$x_{33} = x_{32} - \eta f'(x_{32}) = 0.035$$

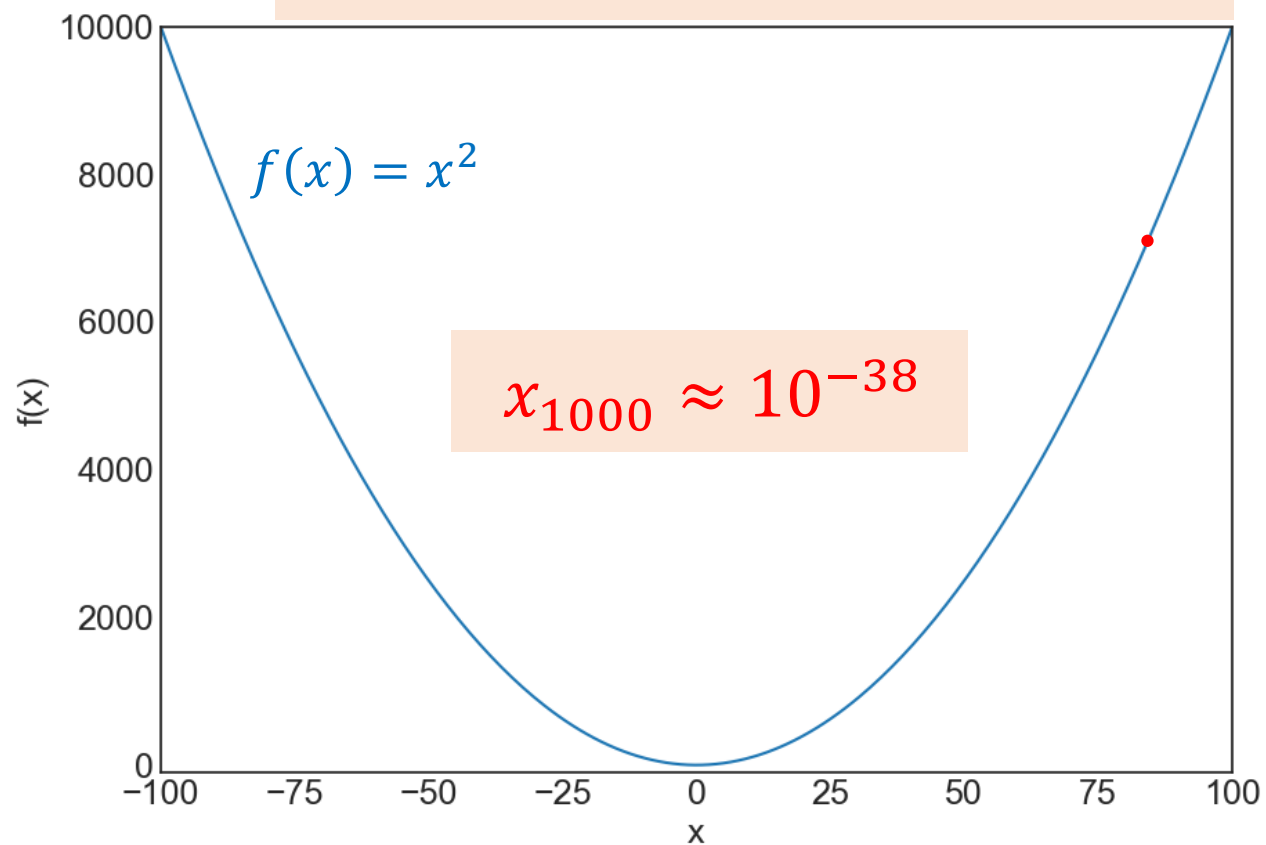
$$f'(x_{34}) = 0.071$$

$$x_{34} = x_{33} - \eta f'(x_{33}) = 0.028$$

Optimization

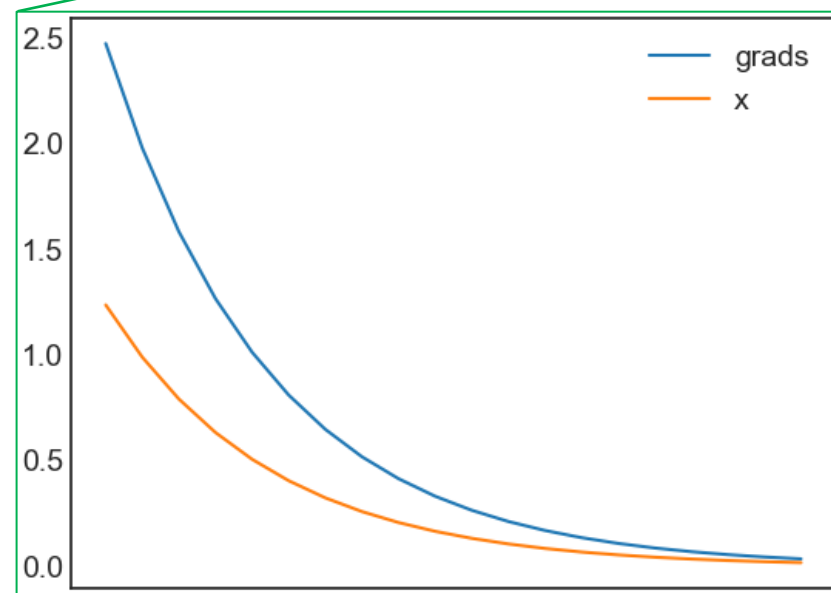
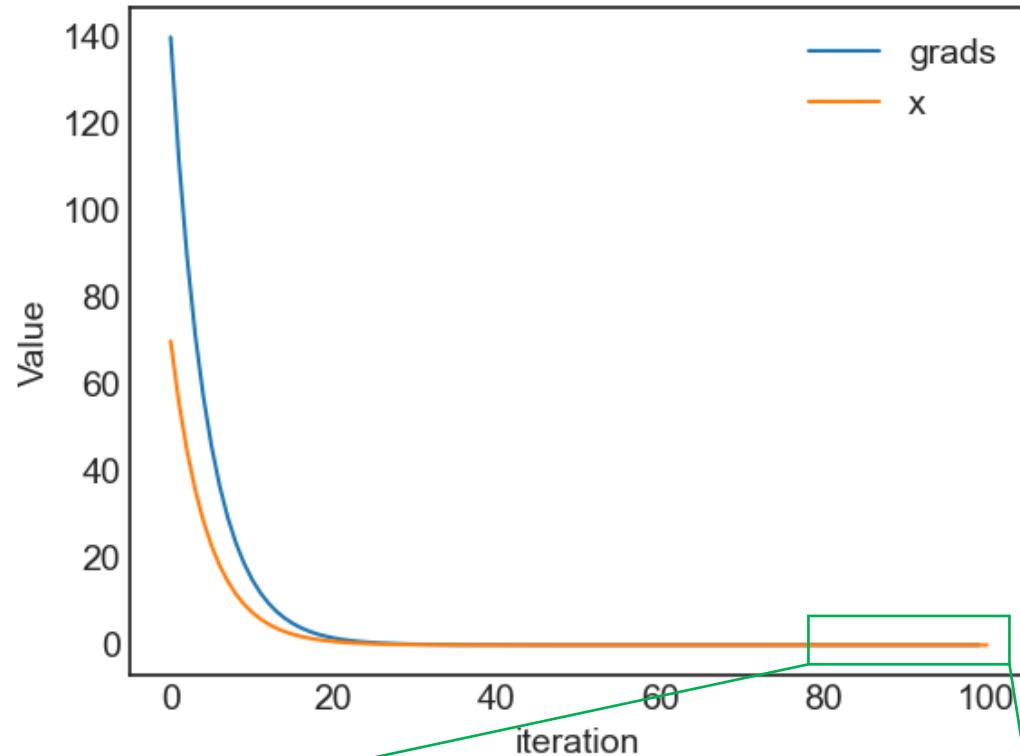
❖ Square function

$$x_t = x_{t-1} - \eta f'(x)$$



$$x_{1000} \approx 10^{-38}$$

Optimized successfully!



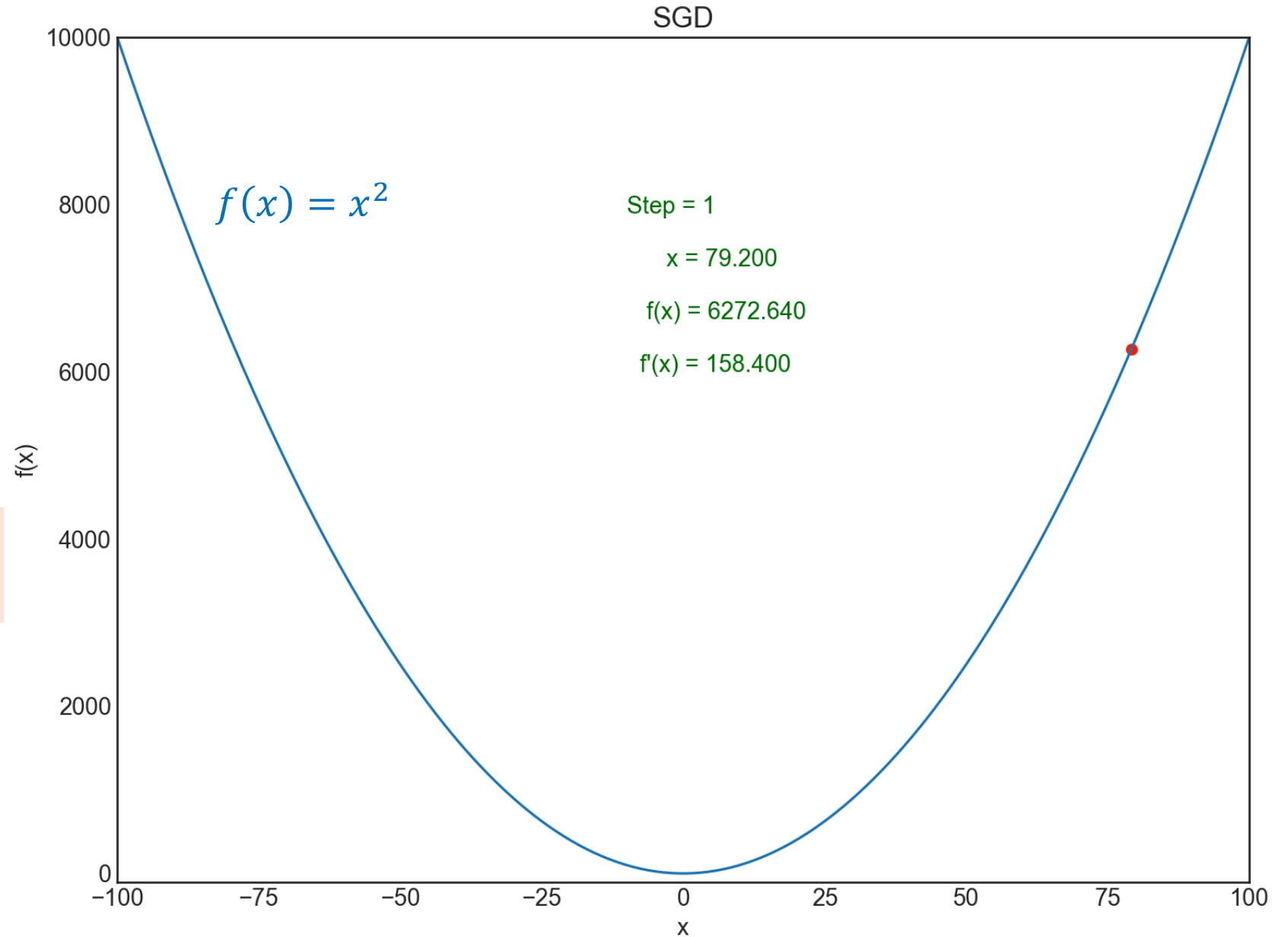
Optimization

❖ Square function

$$x_0 = 99.0$$

$$\eta = 0.1$$

$$x_t = x_{t-1} - \eta f'(x)$$



Optimization

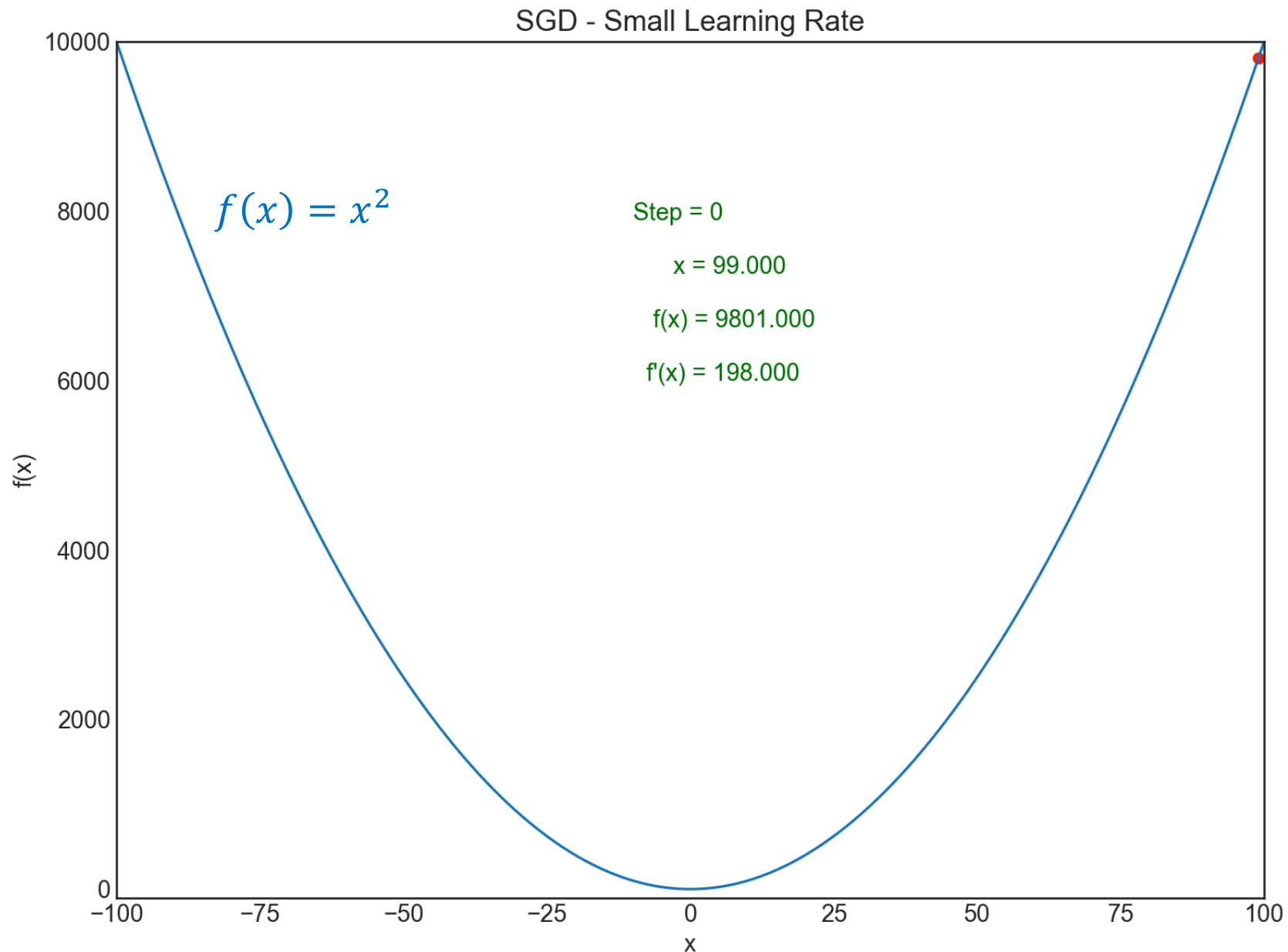
❖ Square function

Discussion

$$x_0 = 99.0$$

$$\eta = 0.001$$

$$x_t = x_{t-1} - \eta f'(x)$$



Optimization

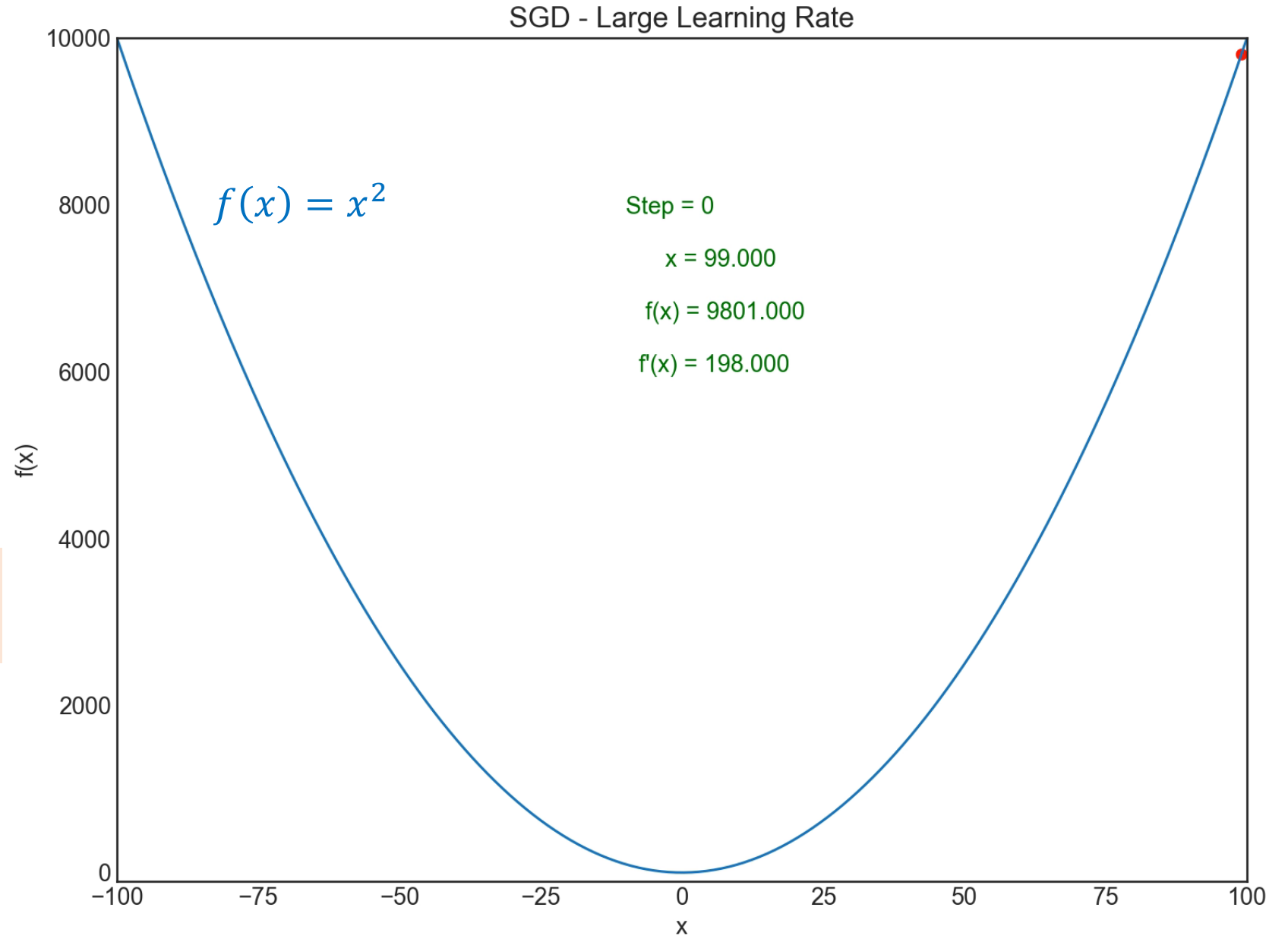
❖ Square function

Discussion

$$x_0 = 99.0$$

$$\eta = 0.8$$

$$x_t = x_{t-1} - \eta f'(x)$$



Optimization

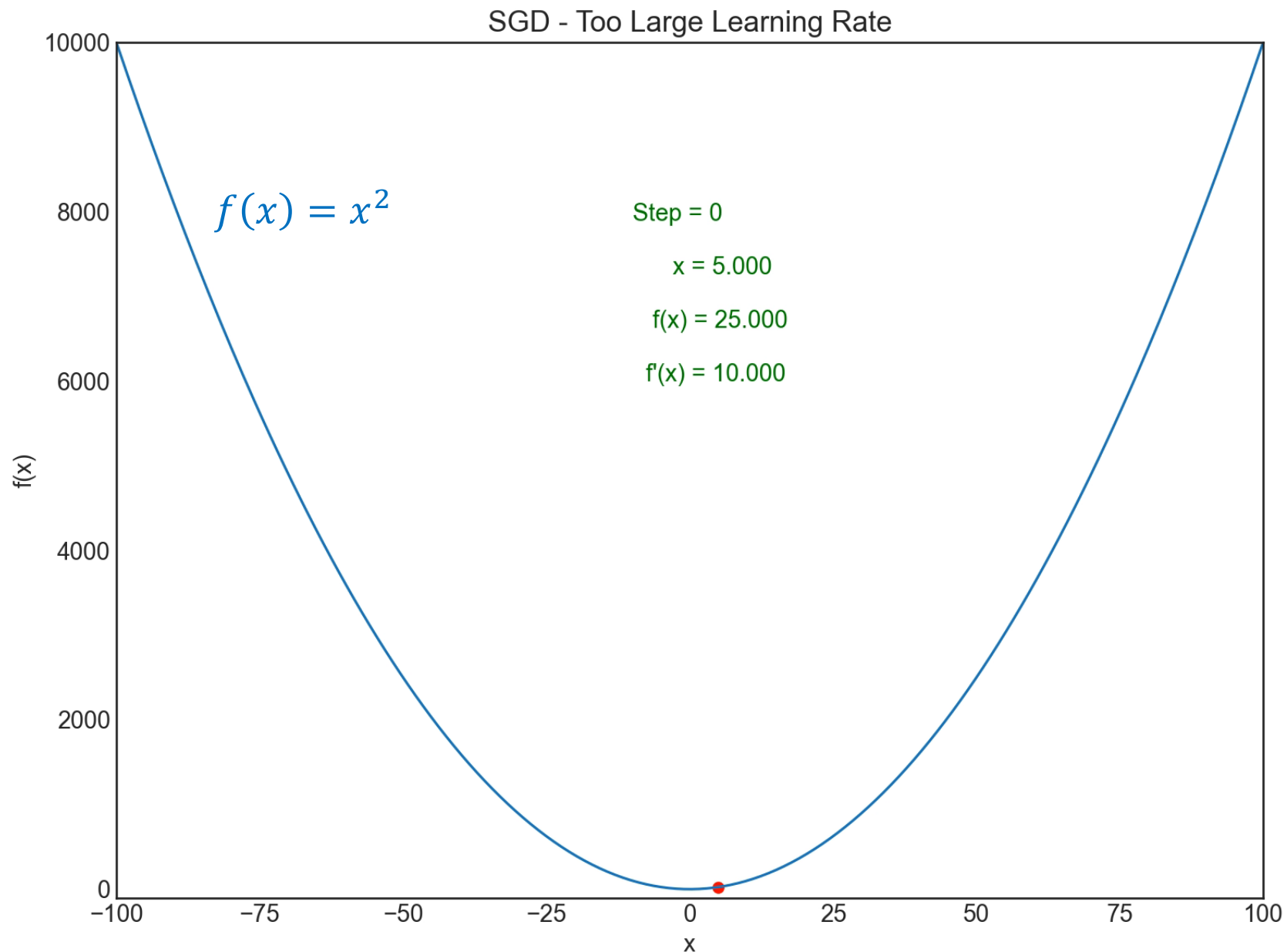
❖ Square function

Discussion

$$x_0 = 99.0$$

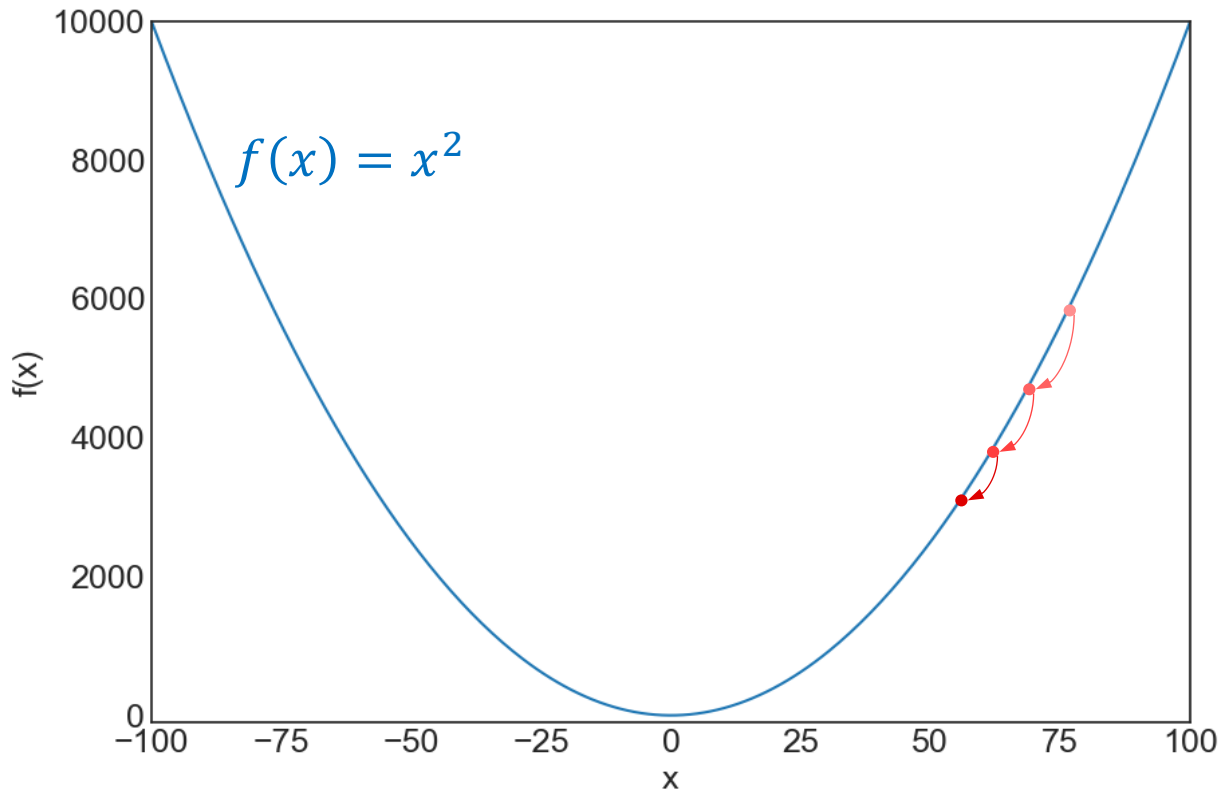
$$\eta = 1.1$$

$$x_t = x_{t-1} - \eta f'(x)$$



Optimization

❖ Square function: Summary



- Given a function $f(x)$, find optimal x_{opt} so that $f(x_{\text{opt}})$ is minimum
- After an update, $f(x_{\text{new}}) \leq f(x_{\text{old}})$

$$x_t = x_{t-1} - \eta f'(x_{t-1})$$

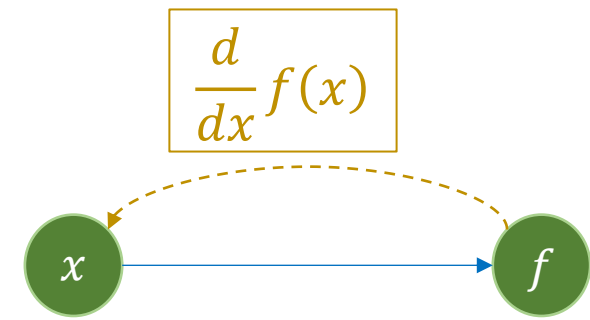
$$x_{30} = 0.069 \quad \eta = 0.1$$

$$f'(x_{30}) = 0.138$$

$$x_{31} = x_{30} - \eta f'(x_{30}) = 0.055$$

$$f'(x_{31}) = 0.11$$

$$x_{32} = x_{31} - \eta f'(x_{31}) = 0.044$$



what does $f(x)$ mean?

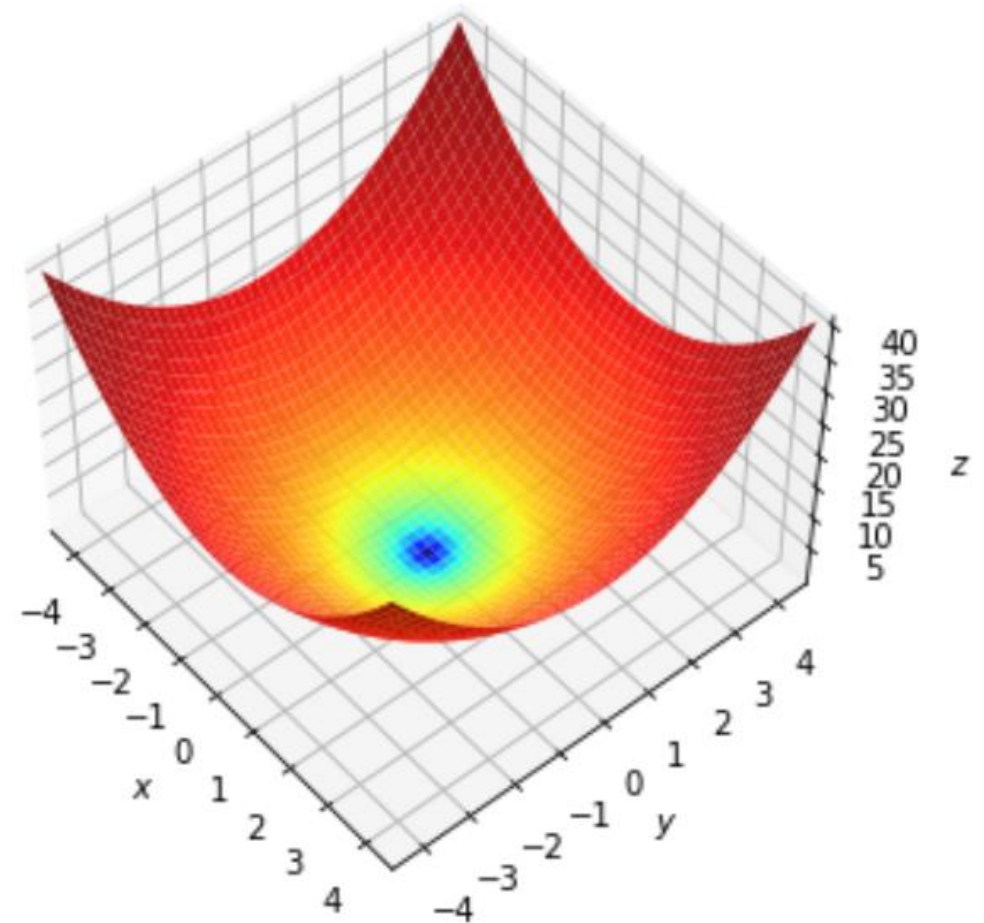
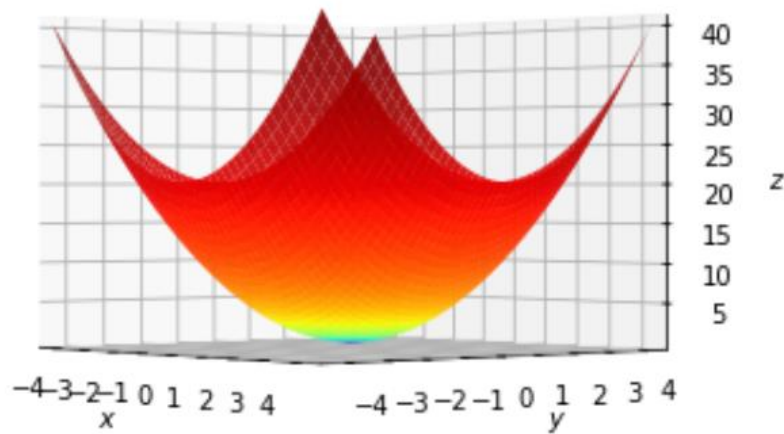
Observation 2

❖ Optimization: 2D function

$$f(x, y) = x^2 + y^2$$

$$-100 \leq x, y \leq 100$$

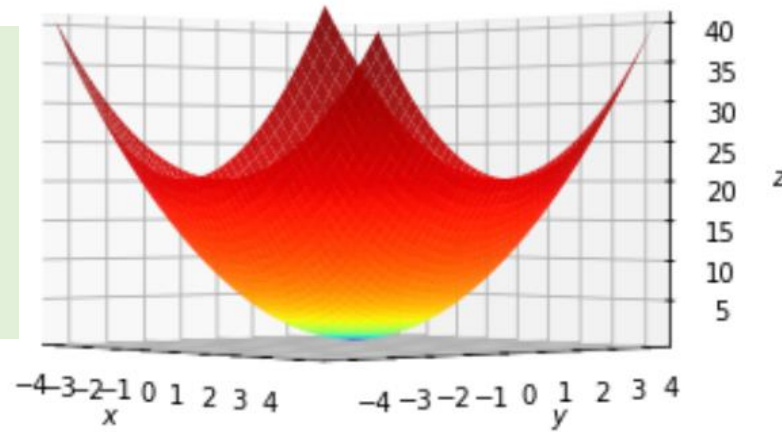
$$x, y \in \mathbb{N}$$



Optimization

❖ Optimization: 2D function

$$f(x, y) = x^2 + y^2$$
$$-100 \leq x, y \leq 100$$
$$x, y \in \mathbb{N}$$



$$x = x - \eta \frac{\partial f(x, y)}{\partial x}$$

$$y = y - \eta \frac{\partial f(x, y)}{\partial y}$$

$$x_0 = 6.0 \quad y_0 = 9.0 \quad \eta = 0.1$$

$$\frac{\partial f(x_0, y_0)}{\partial x} = 12 \quad \frac{\partial f(x_0, y_0)}{\partial y} = 18$$
$$x_1 = 4.8 \quad y_1 = 7.2$$

$$\frac{\partial f(x_1, y_1)}{\partial x} = 9.6 \quad \frac{\partial f(x_1, y_1)}{\partial y} = 14.4$$
$$x_2 = 3.84 \quad y_2 = 5.75$$

$$\frac{\partial f(x_2, y_2)}{\partial x} = 7.68 \quad \frac{\partial f(x_2, y_2)}{\partial y} = 11.51$$
$$x_3 = 3.07 \quad y_3 = 4.608$$

$$\frac{\partial f(x_3, y_3)}{\partial x} = 6.14 \quad \frac{\partial f(x_3, y_3)}{\partial y} = 9.21$$
$$x_4 = 2.45 \quad y_4 = 3.68$$

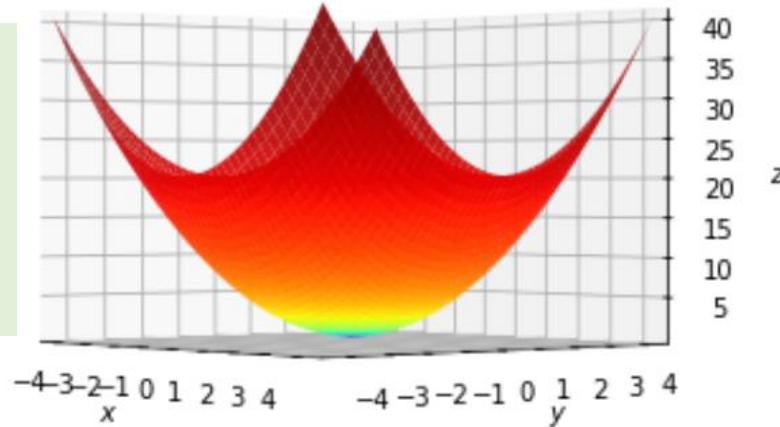
Optimization

❖ Optimization: 2D function

$$f(x, y) = x^2 + y^2$$

$$-100 \leq x, y \leq 100$$

$$x, y \in \mathbb{N}$$

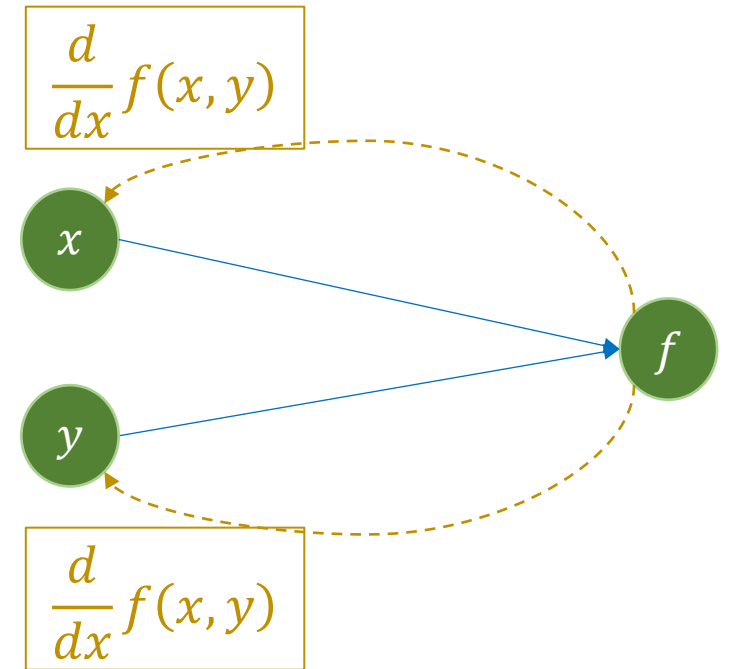


Summary:

- Given a function $f(x, y)$, find optimal $(x_{\text{opt}}, y_{\text{opt}})$ so that $f(x_{\text{opt}}, y_{\text{opt}})$ is minimum
- After an update, $f(x_{\text{new}}, y_{\text{new}}) \leq f(x_{\text{old}}, y_{\text{old}})$

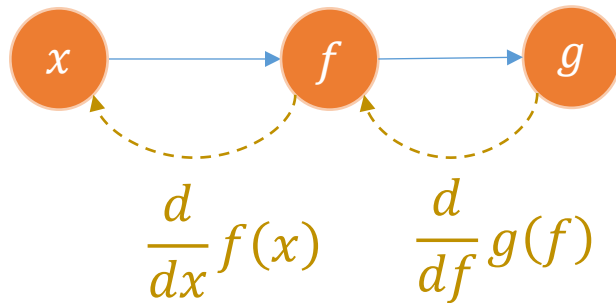
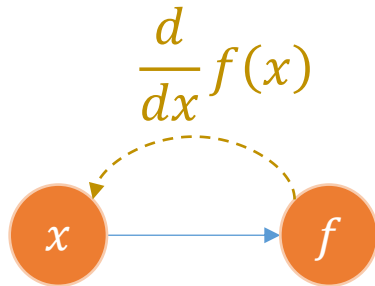
what does $f(x, y)$ mean?

$$\begin{aligned} \frac{\partial f(x_2, y_2)}{\partial x} &= 7.68 & \frac{\partial f(x_2, y_2)}{\partial y} &= 11.51 \\ x_3 &= 3.07 & y_3 &= 4.608 \end{aligned}$$
$$\begin{aligned} \frac{\partial f(x_3, y_3)}{\partial x} &= 6.14 & \frac{\partial f(x_3, y_3)}{\partial y} &= 9.21 \\ x_4 &= 2.45 & y_4 &= 3.68 \end{aligned}$$

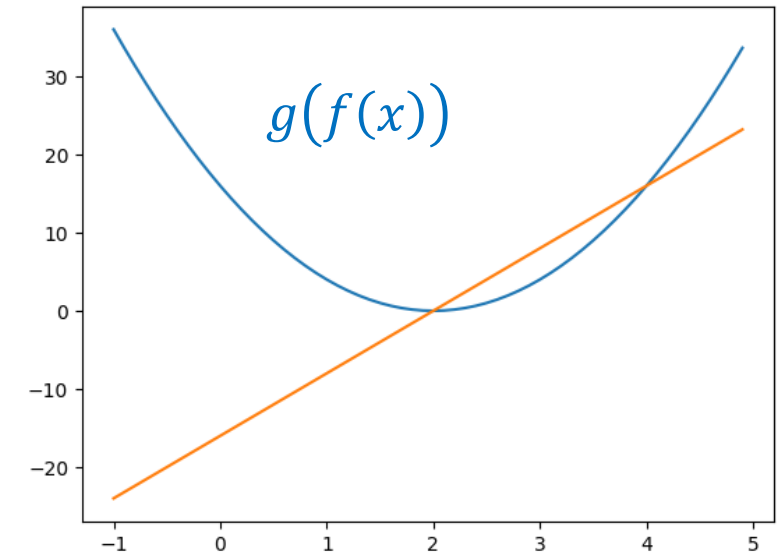
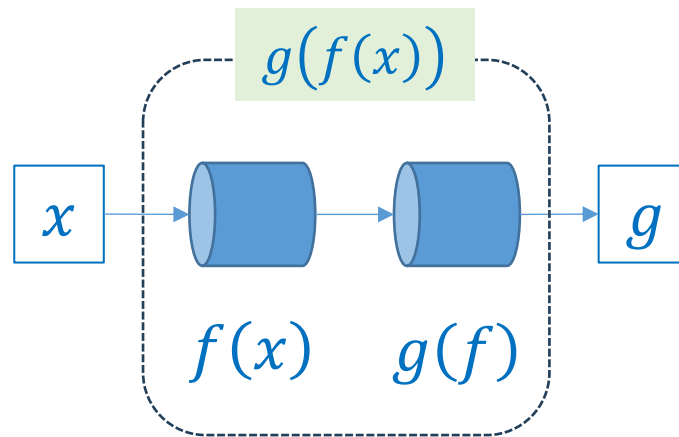
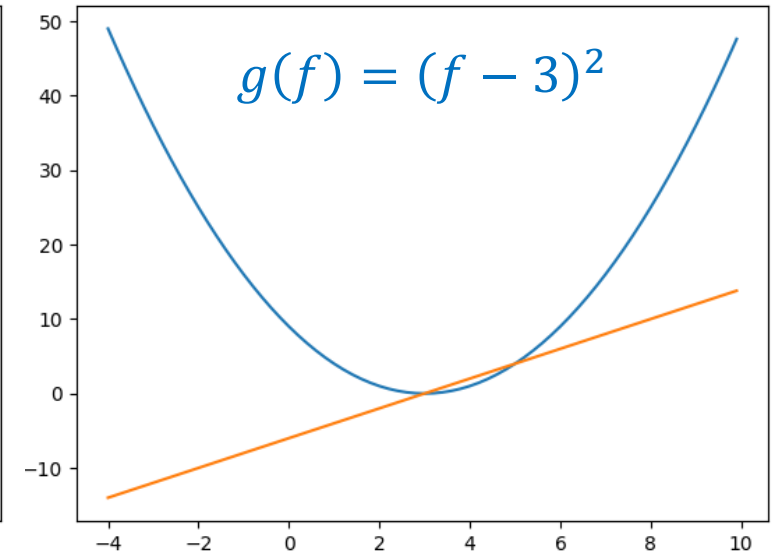
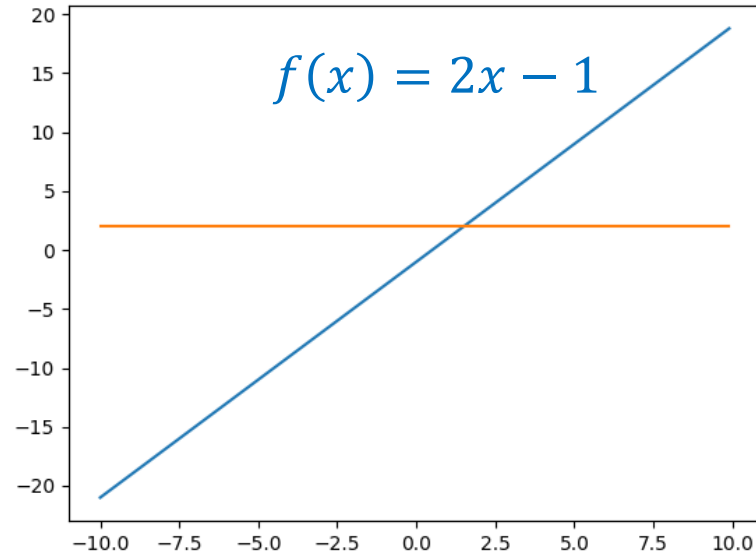


Observation 3

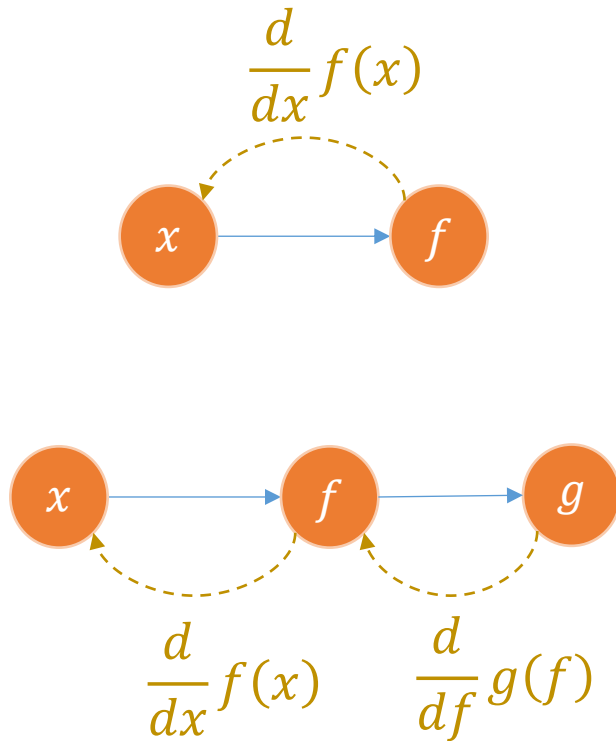
❖ For composite function



$$\frac{d}{dx} g(f(x)) = \left[\frac{d}{df} g(f) \right] * \left[\frac{d}{dx} f(x) \right]$$



❖ For composite function and chain rule



$$\frac{d}{dx} g(f(x)) = \left[\frac{d}{df} g(f) \right] * \left[\frac{d}{dx} f(x) \right]$$



$$f(x) = 2x - 1$$

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

$$f'(x) = 2$$

$$g'(f) = 2(f - 3)$$



$$\frac{dg}{dx} = \frac{dg}{df} \frac{df}{dx}$$

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

$$= 4(2x - 1 - 3)$$

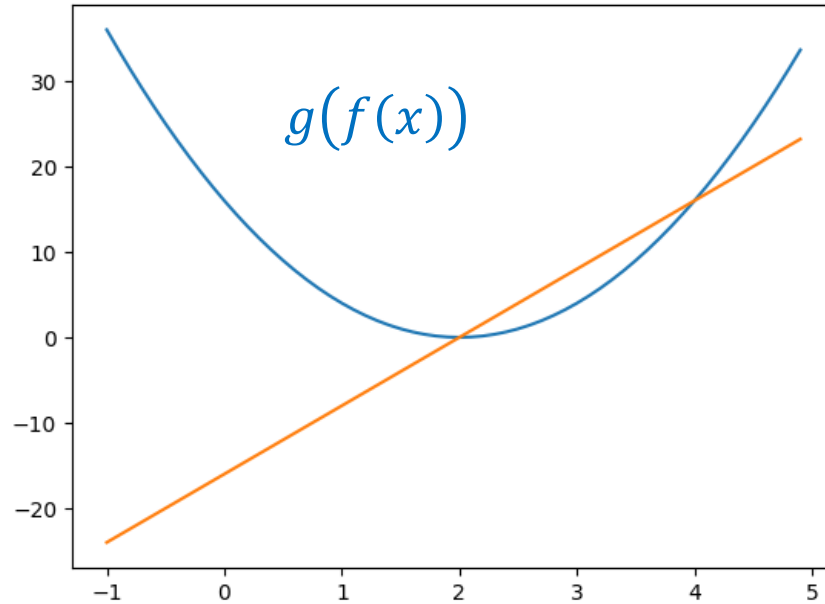
$$= 8x - 16$$

Implementation

$$f(x) = 2x - 1$$

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

$$\frac{dg}{dx} = 8x - 16$$



```
1 def fx(x):  
2     return 2*x - 1  
3  
4 def gf(f):  
5     return (f-3)**2  
6  
7 def dg_dx(x):  
8     return 8*x - 16
```

```
1 import random  
2  
3 # parameters  
4 lr = 0.1  
5  
6 # initialize x  
7 x = 60  
8  
9 old_loss = gf(fx(x)) # Logging  
10 print(f'old_loss: {old_loss}')
```

```
11  
12 # compute derivative  
13 dg_dx_value = dg_dx(x)  
14  
15 # update  
16 x = x - lr*dg_dx_value  
17  
18 new_loss = gf(fx(x)) # Logging  
19 print(f'new_loss: {new_loss}')
```

```
old_loss: 13456
```

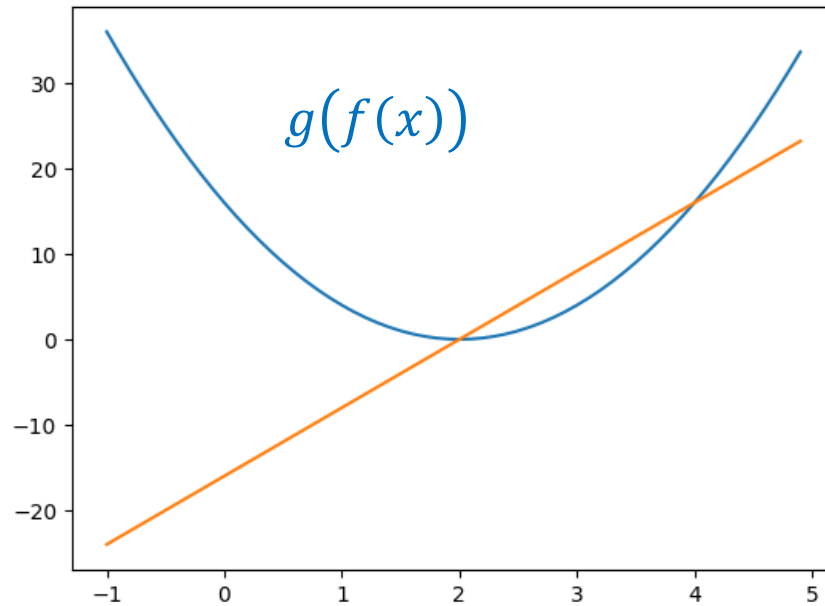
```
new_loss: 538.23999999999994
```

Implementation

$$f(x) = 2x - 1$$

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

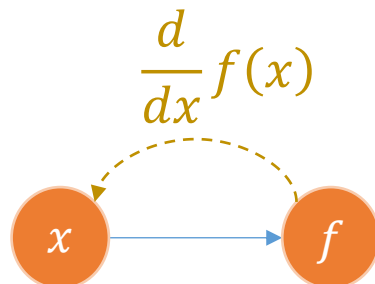
$$\frac{dg}{dx} = 8x - 16$$



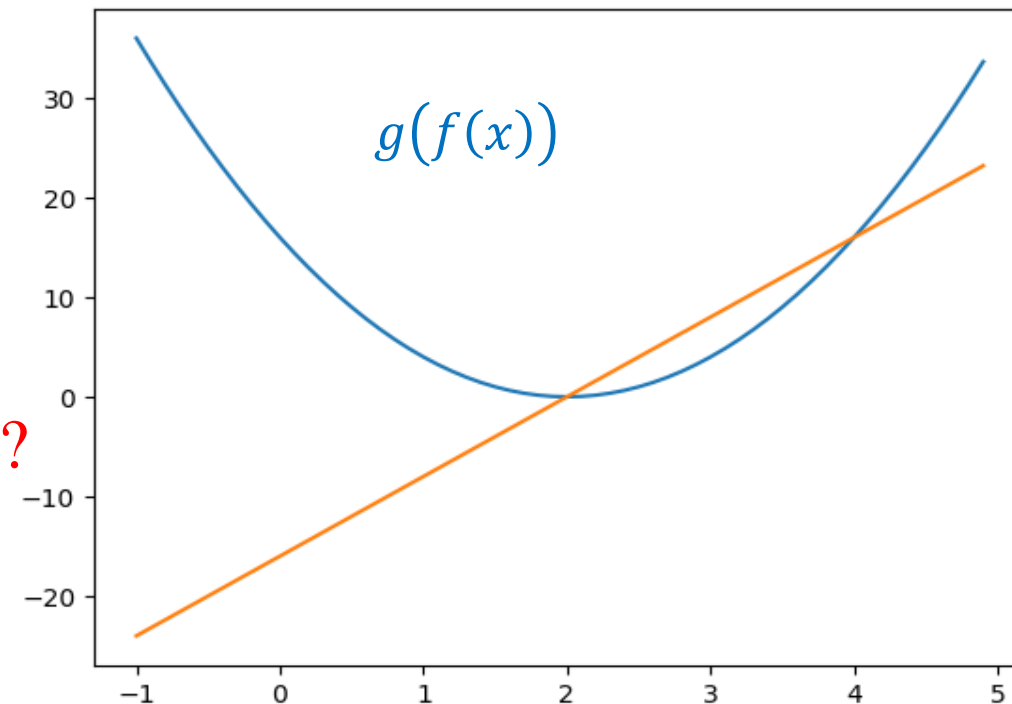
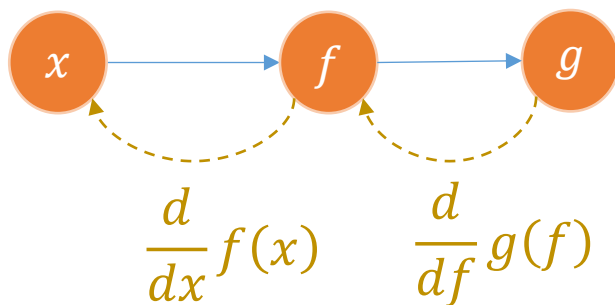
```
1 def fx(x):  
2     return 2*x - 1  
3  
4 def gf(f):  
5     return (f-3)**2  
6  
7 def dg_dx(x):  
8     return 8*x - 16
```

```
1 import random  
2  
3 # parameters  
4 num_steps = 5  
5 lr = 0.1  
6  
7 # set x randomly  
8 x = random.randint(-100, 100)  
9  
10 for _ in range(num_steps):  
11     # logging  
12     loss = gf(fx(x))  
13  
14     # compute derivative  
15     dg_dx_value = dg_dx(x)  
16  
17     # update  
18     x = x - lr*dg_dx_value
```

❖ For composite function: Summary



what does $g(f(x))$ mean?



- Given a function $g(f(x))$, find optimal x_{opt} so that $f(x_{\text{opt}})$ is minimum
- After an update, $g(f(x_{\text{new}})) \leq g(f(x_{\text{old}}))$

$$\frac{d}{dx} g(f(x)) = \left[\frac{d}{df} g(f) \right] * \left[\frac{d}{dx} f(x) \right]$$

Outline

SECTION 1

Optimization

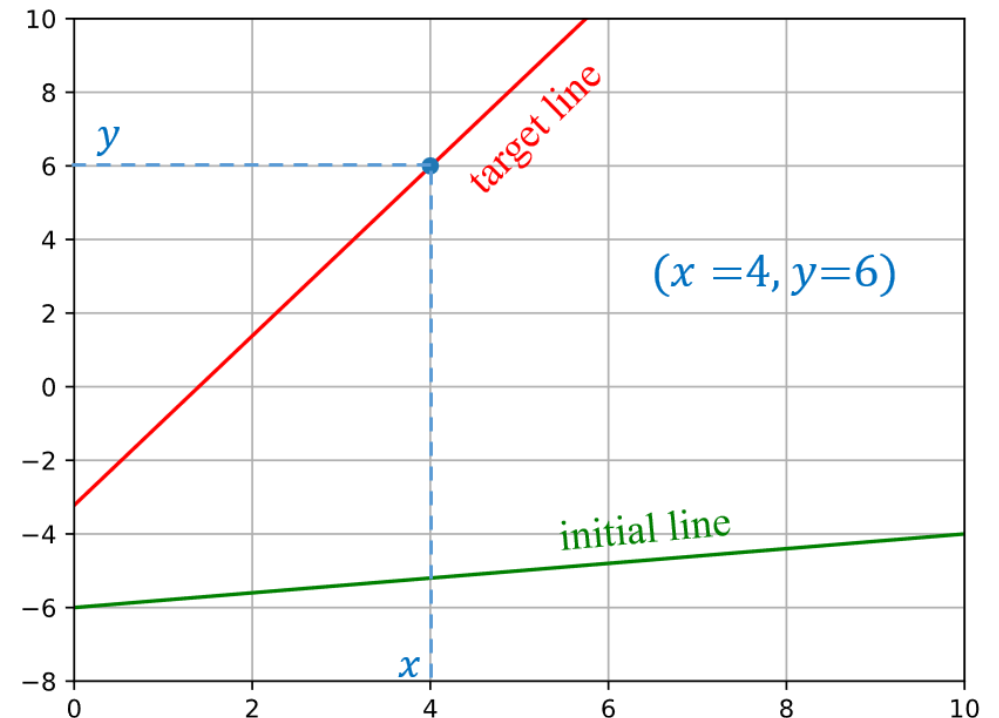
SECTION 2

Problem Solving

SECTION 3

Towards Linear Regression

$$f(x) = ax + b$$

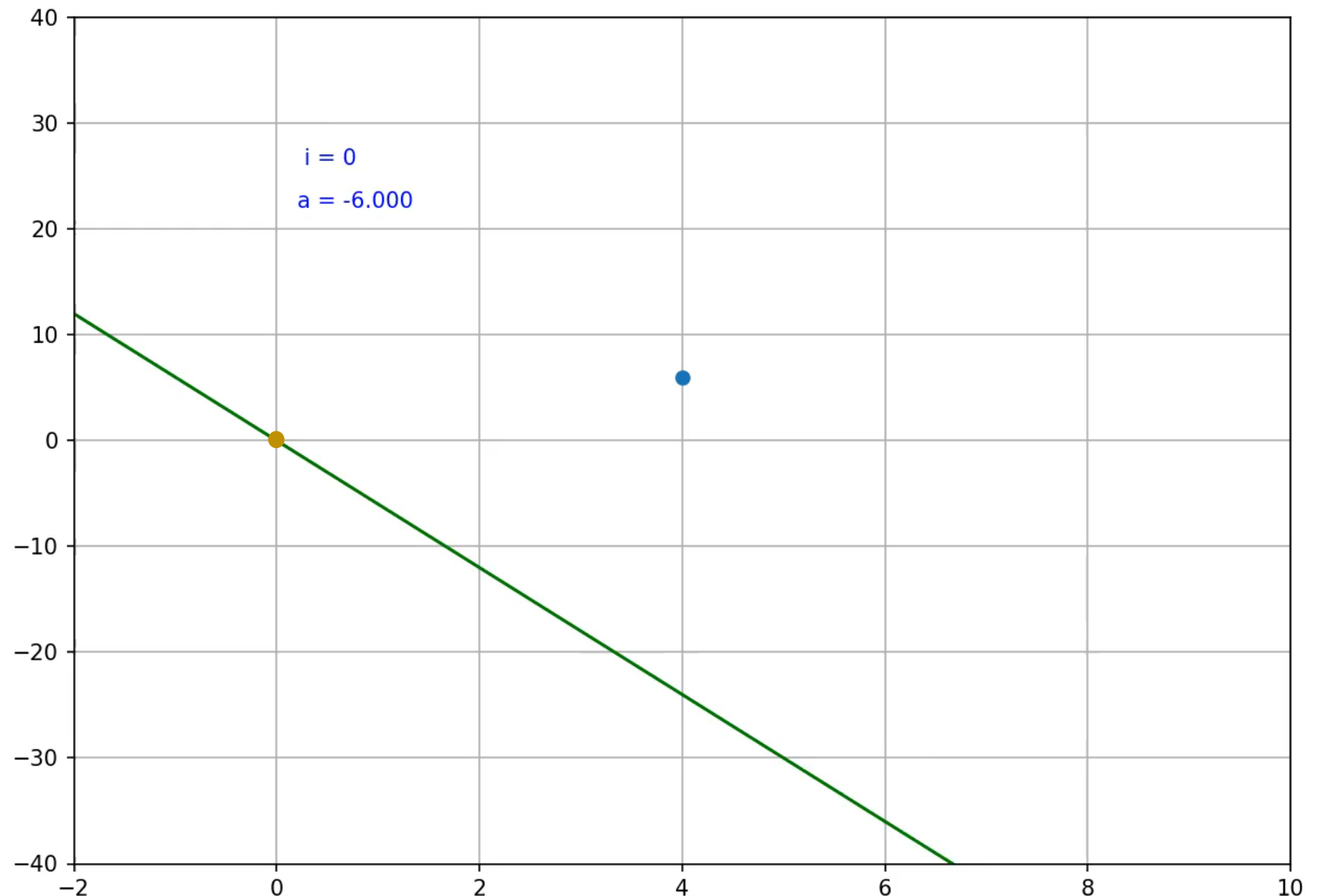


Gradient-based Optimization

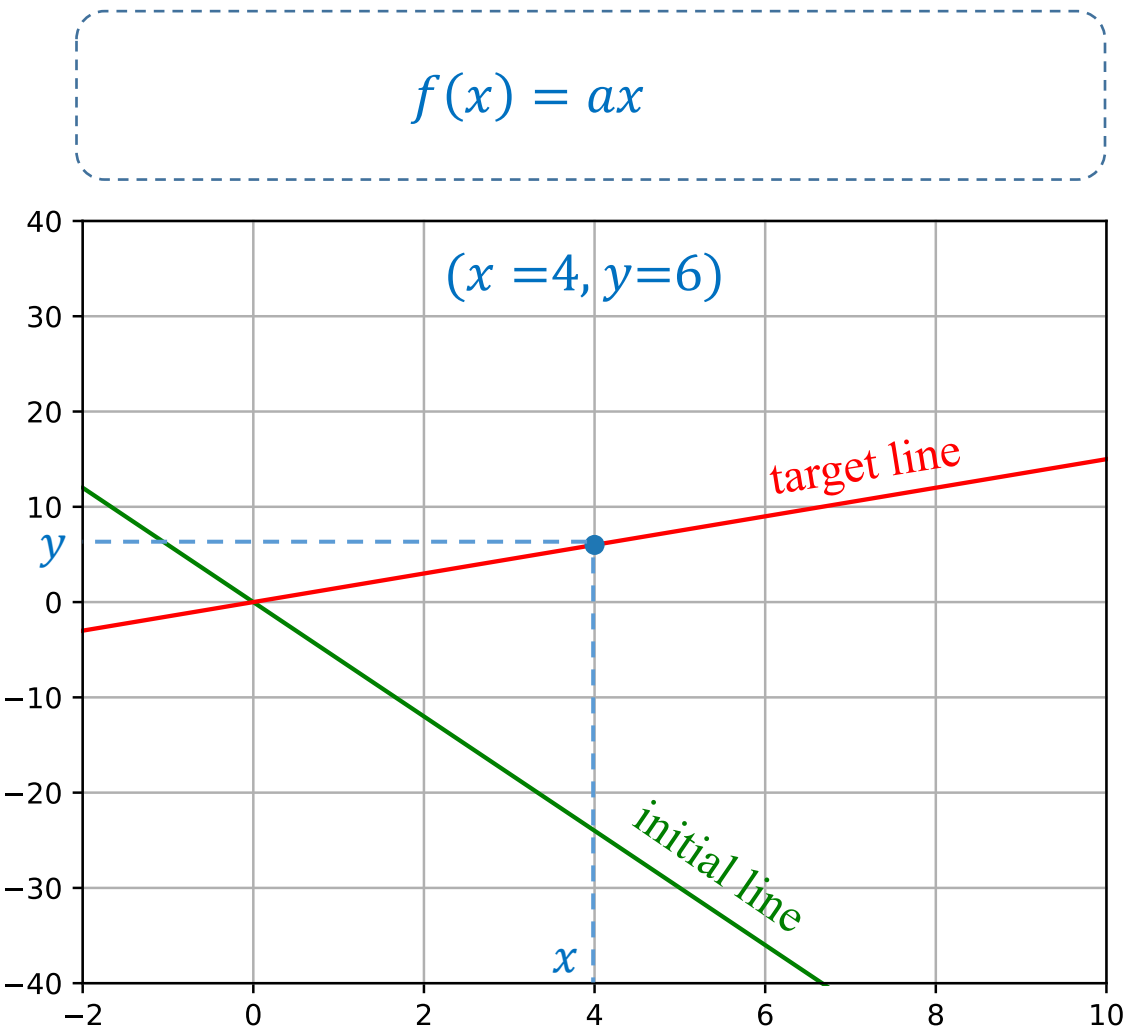
❖ Problem 1

Given a line that is always
through the origin

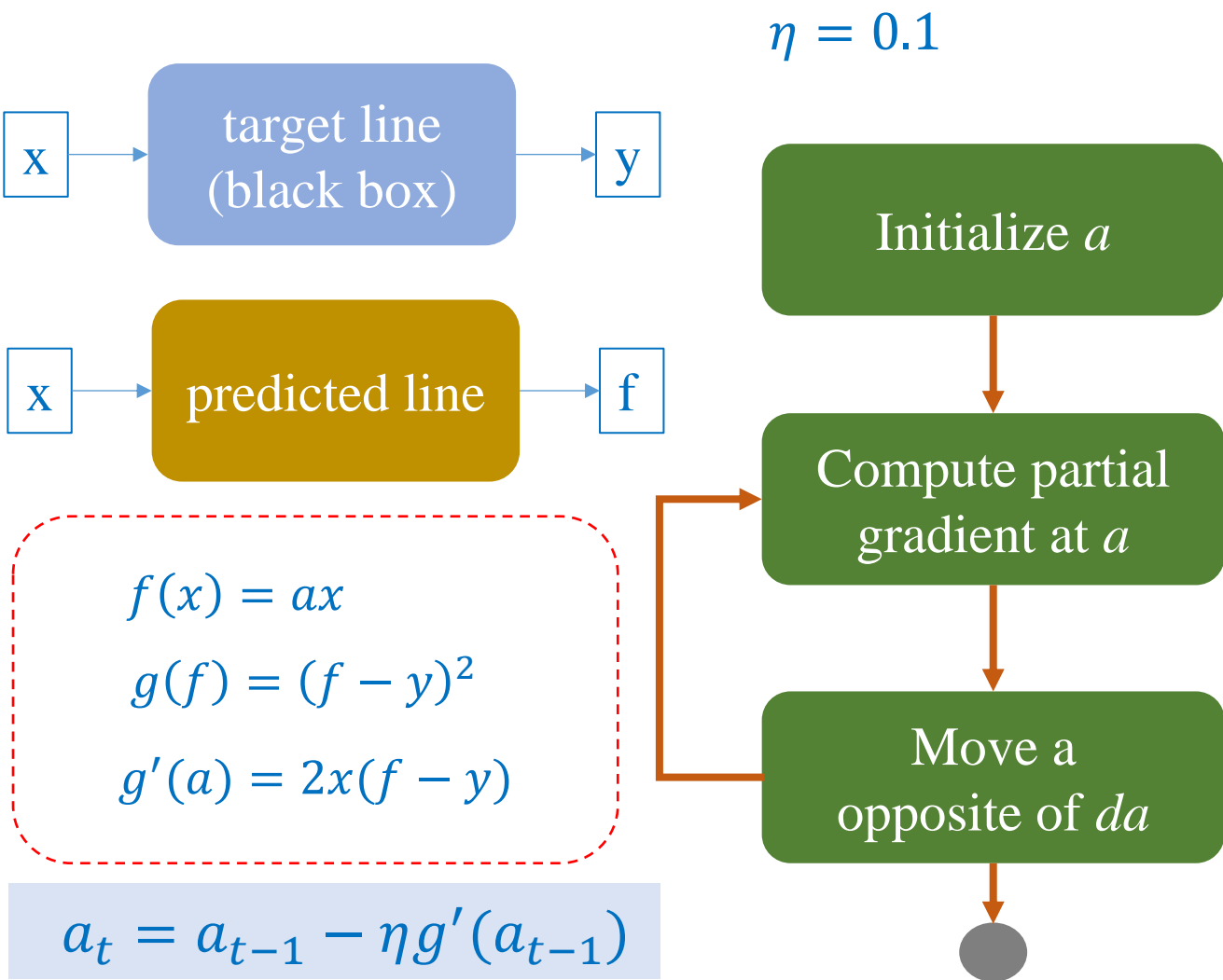
How to move the green line
so that it is also through the
blue point?



❖ Problem 1



❖ Constraints



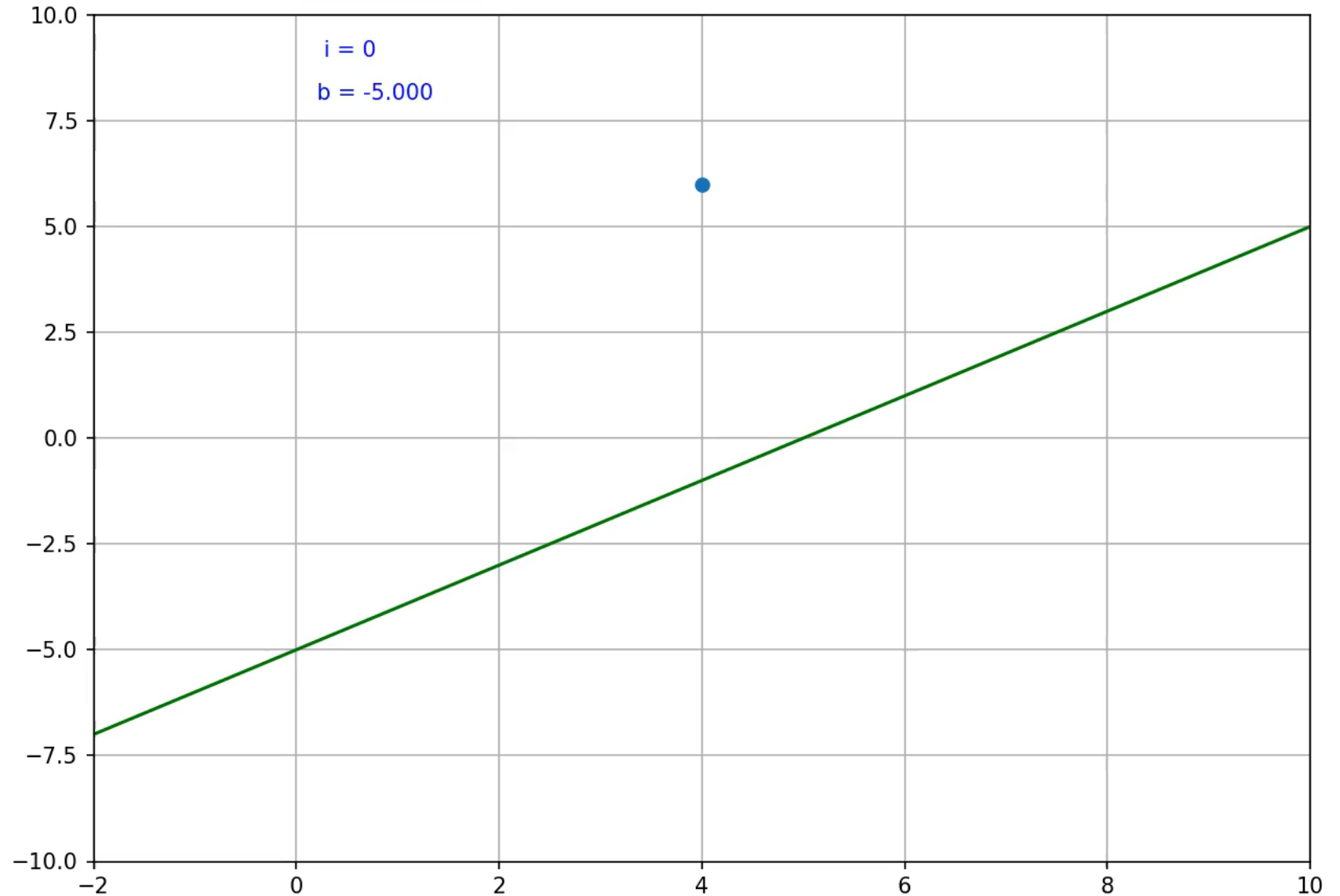
Gradient-based Optimization

❖ Problem 2

Given a line

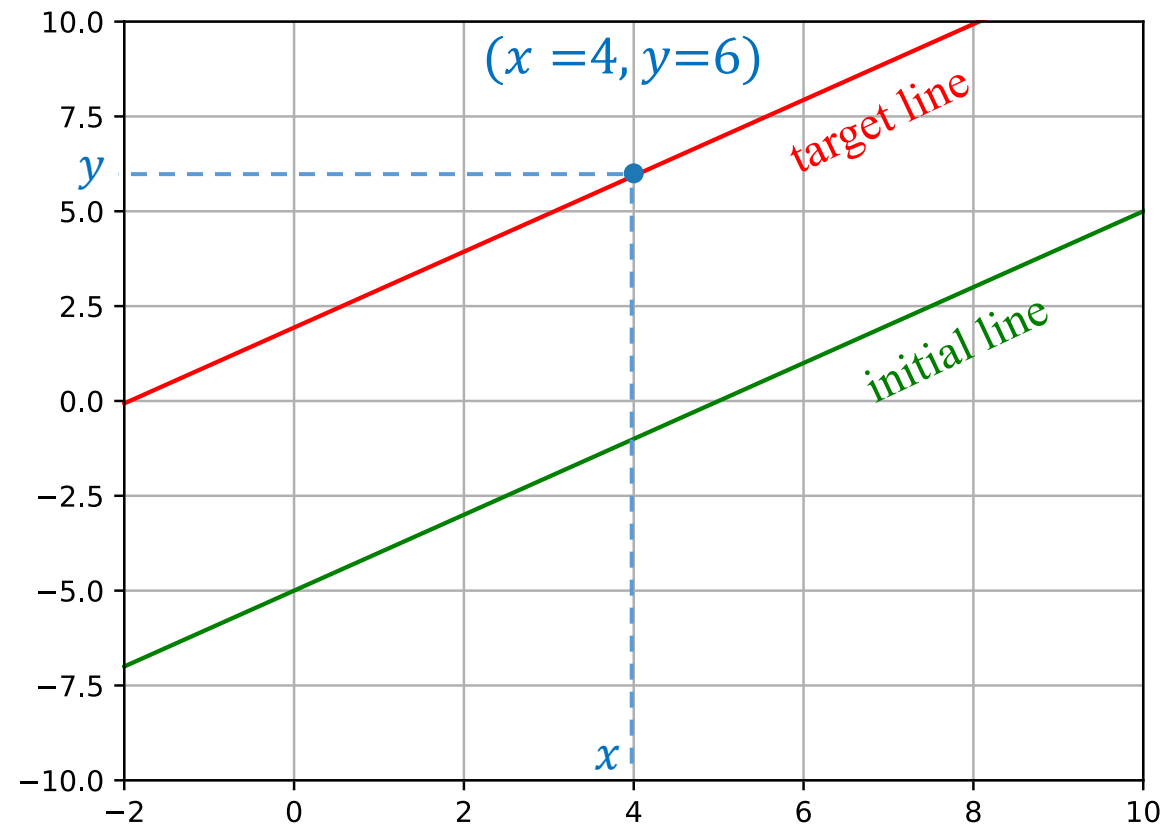
$$f(x) = x + b$$

How to move the green line
so that it is also through the
blue point?



❖ Problem 2

$$f(x) = x + b$$



❖ Constraints



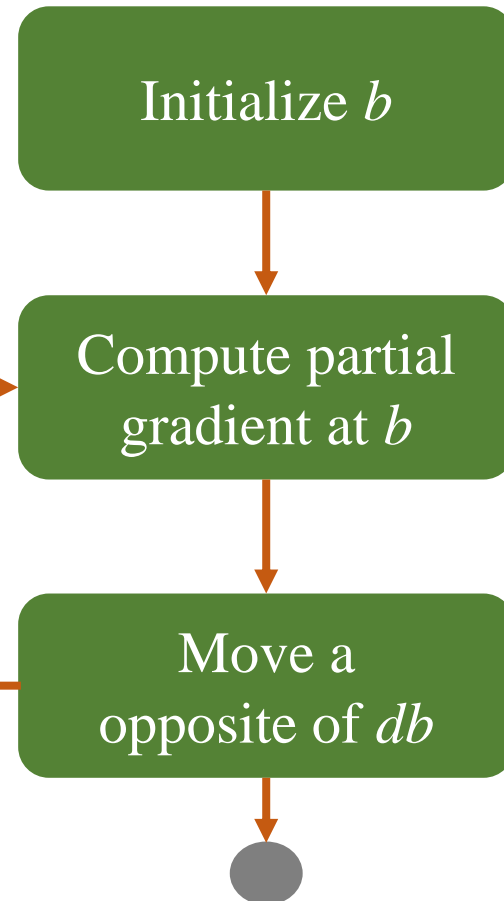
$$f(x) = x + b$$

$$g(f) = (f - y)^2$$

$$g'(b) = 2(f - y)$$

$$b_t = b_{t-1} - \eta g'(b_{t-1})$$

$$\eta = 0.1$$



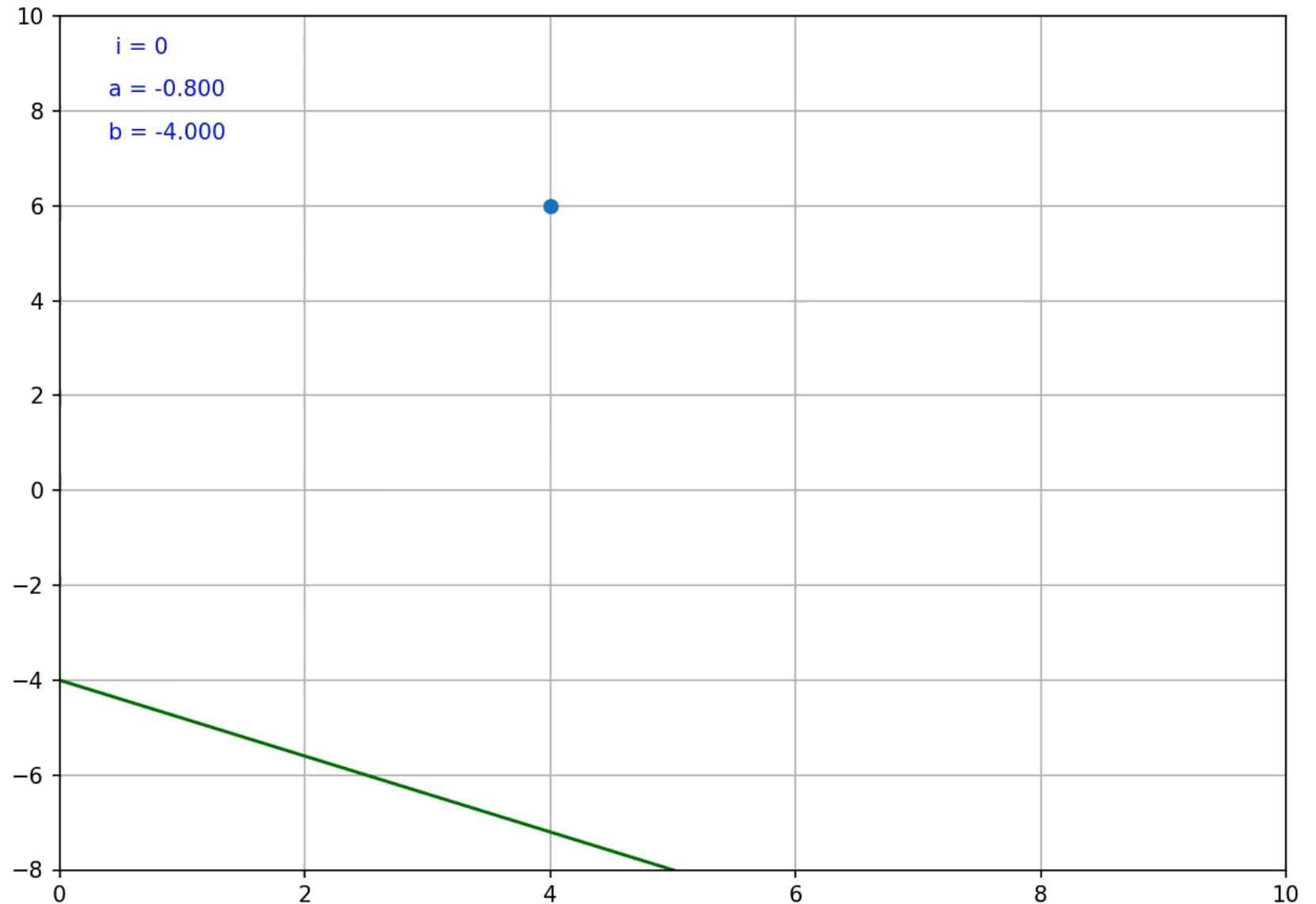
Gradient-based Optimization

❖ Problem 3

Given a line

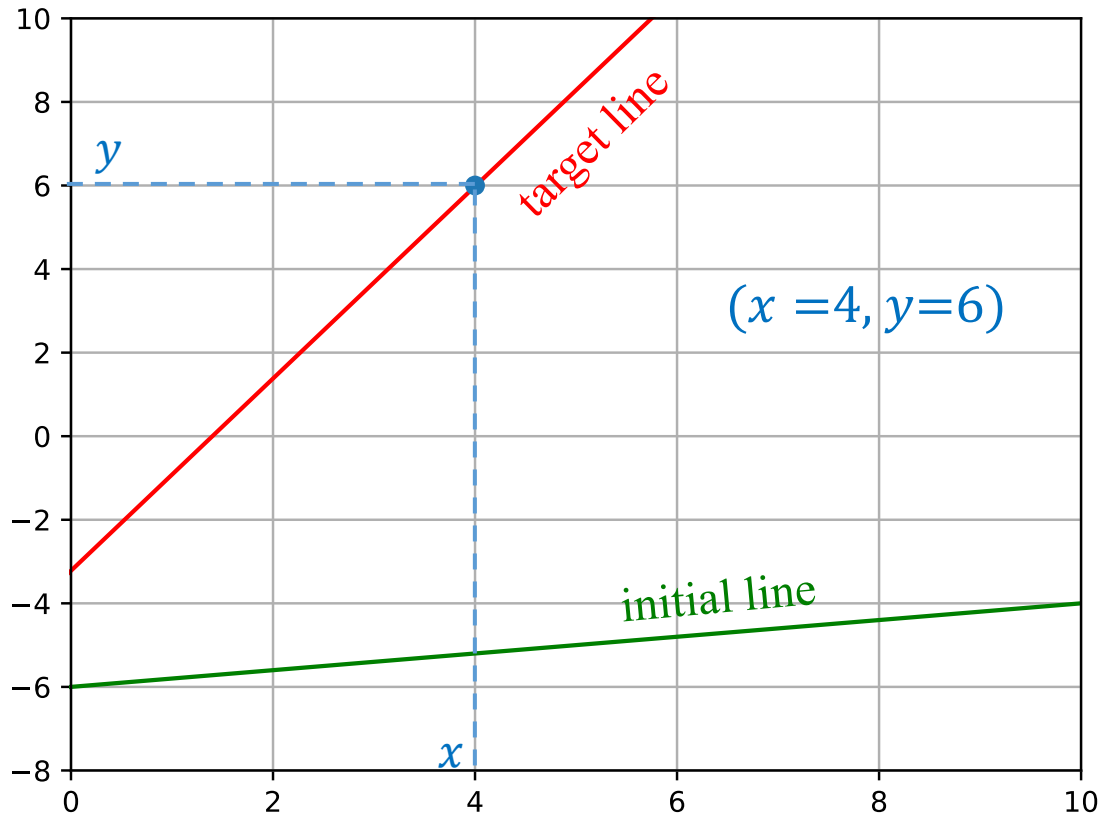
$$f(x) = ax + b$$

How to move the green line
so that it is also through the
blue point?



❖ Problem 3

$$f(x) = ax + b$$



❖ Constraints



$$f(x) = ax + b$$

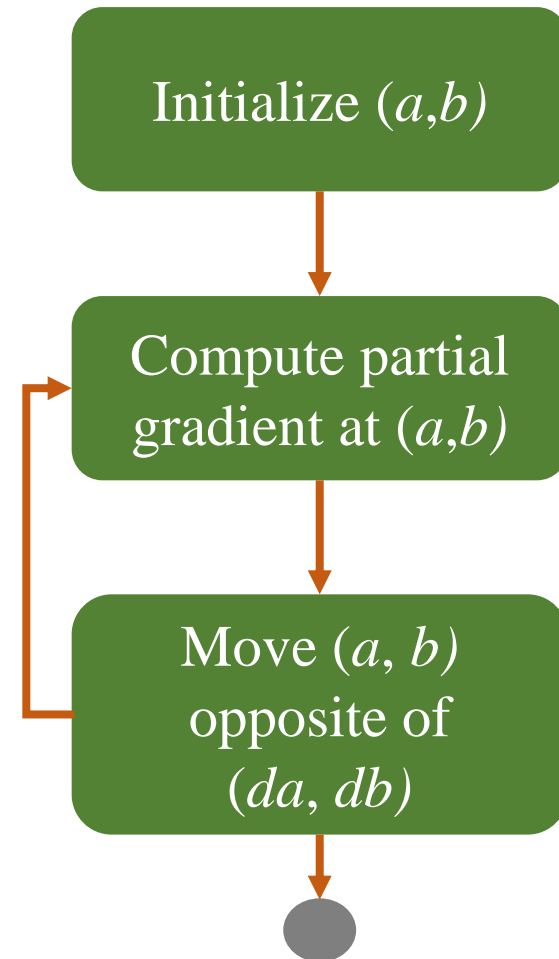
$$g(f) = (f - y)^2$$

$$g'(a) = 2x(f - y)$$

$$g'(b) = 2(f - y)$$

$$\theta_t = \theta_{t-1} - \eta g'(\theta_{t-1})$$

$$\eta = 0.1$$



Outline

SECTION 1

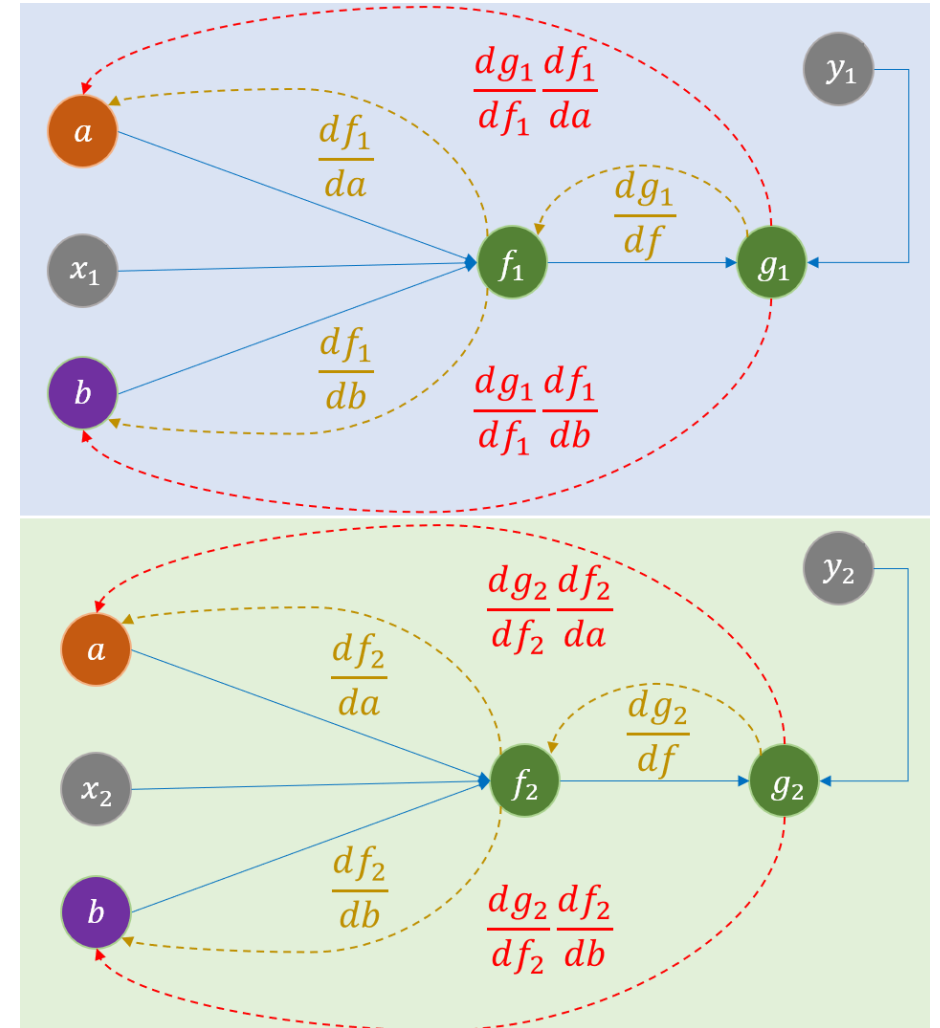
Optimization

SECTION 2

Problem Solving

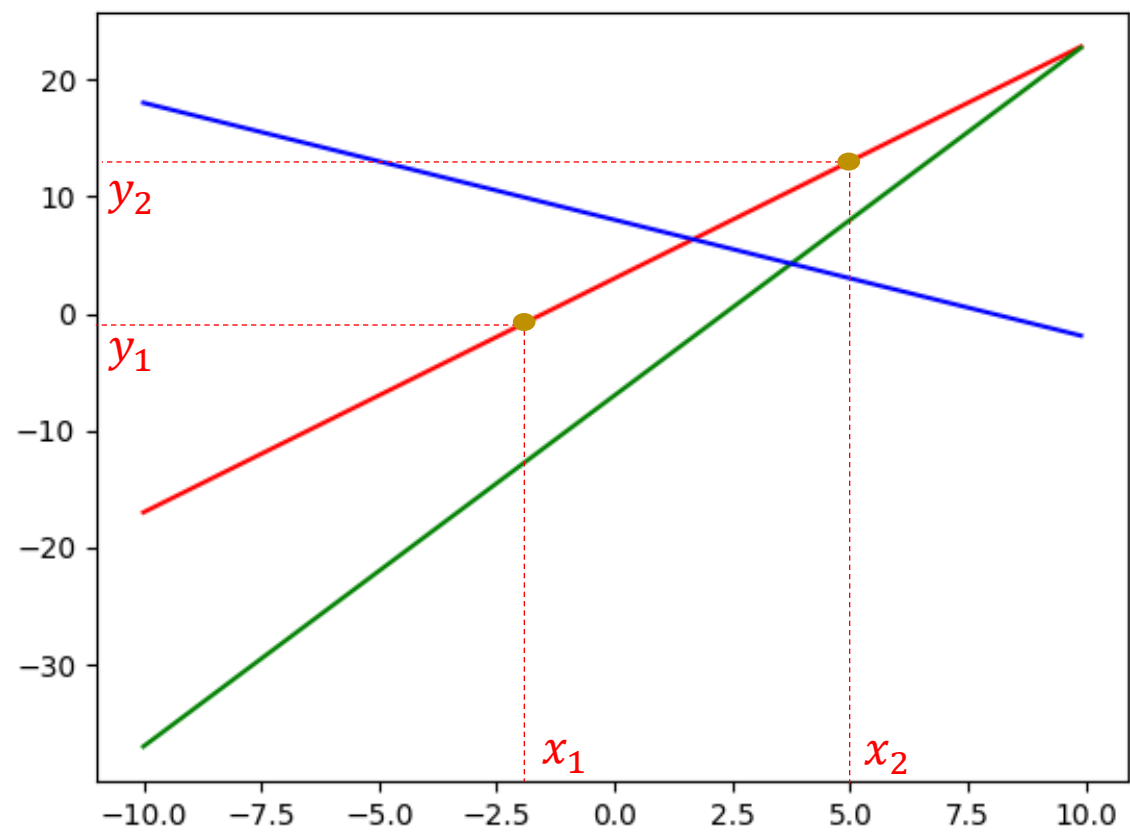
SECTION 3

Towards Linear Regression



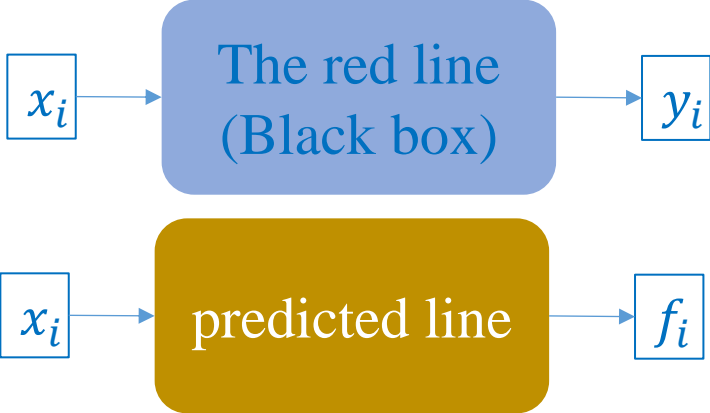
❖ What about having two samples

$f(x) = ax + b$

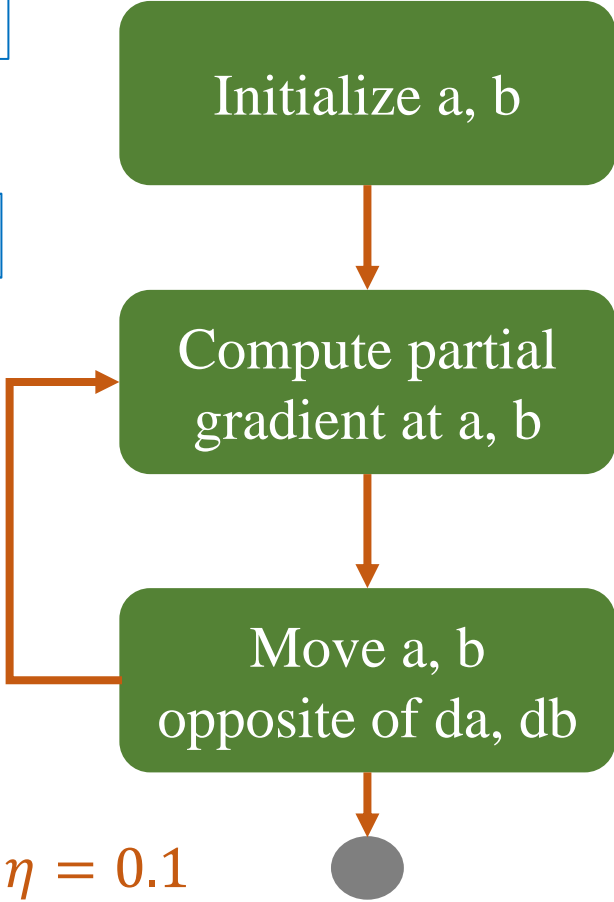


❖ Constraints

$$\theta_t = \theta_{t-1} - \eta f'(\theta_{t-1})$$



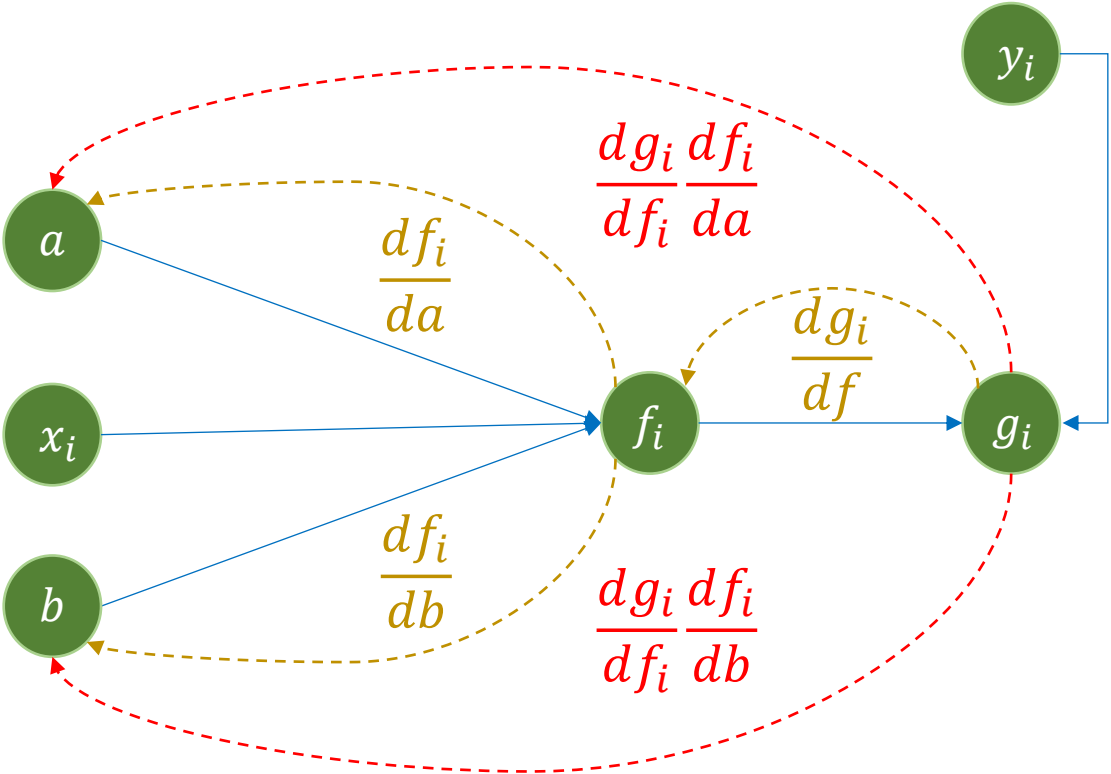
$$(x_1 = -2, y_1 = -1)$$
$$(x_2 = 5, y_2 = 13)$$
$$g(f) = (f - y)^2$$



❖ Equations for partial gradients

$$f(x_i) = ax_i + b \qquad (x_1=1, y_1=5)$$

$$g(f_i) = (f_i - y_i)^2 \qquad (x_2=2, y_2=7)$$



$$\frac{df}{da} = x \qquad \frac{df}{db} = 1$$

$$\frac{dg}{df} = 2(f - y)$$

$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$

$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$

During looking for optimal a and b, at a given time, a and b have concrete values

❖ Optimization for a composite function

Find a and b so that $g(f(x))$ is minimum

$$f(x_i) = ax_i + b \quad (x_1=1, y_1=5)$$

$$g(f_i) = (f_i - y_i)^2 \quad (x_2=2, y_2=7)$$

Partial derivative functions

$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$

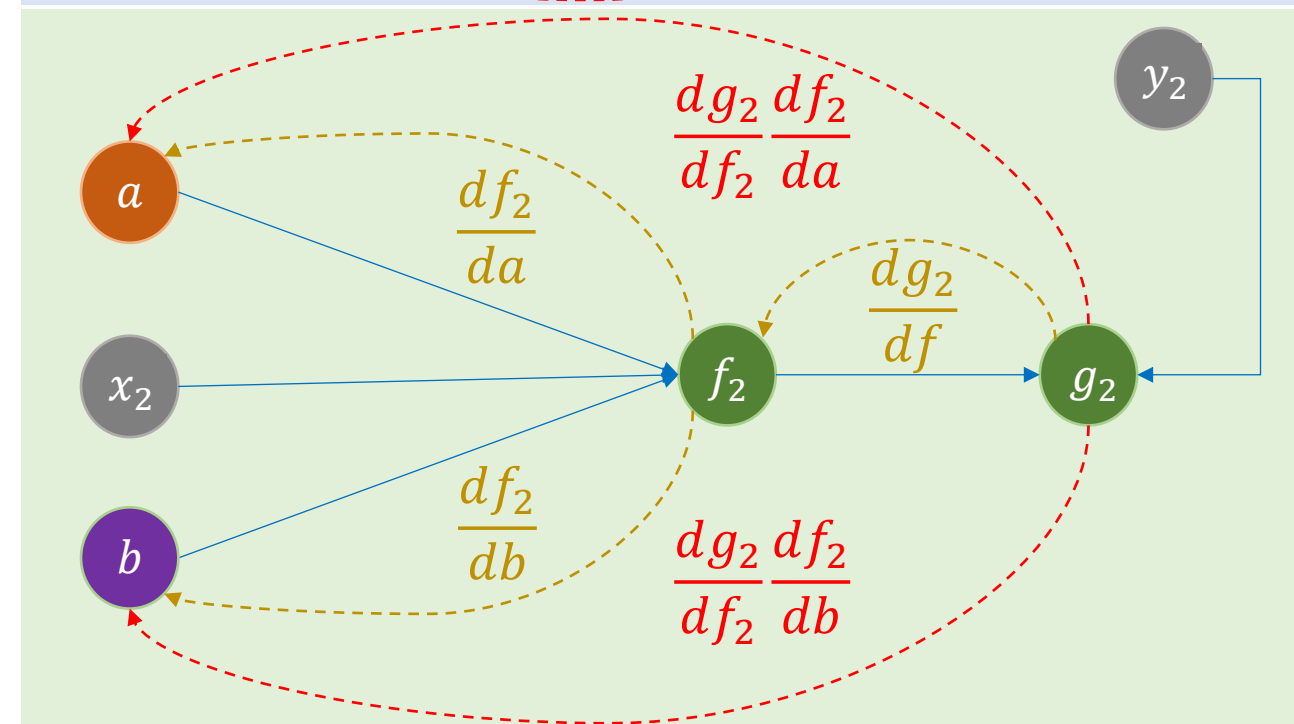
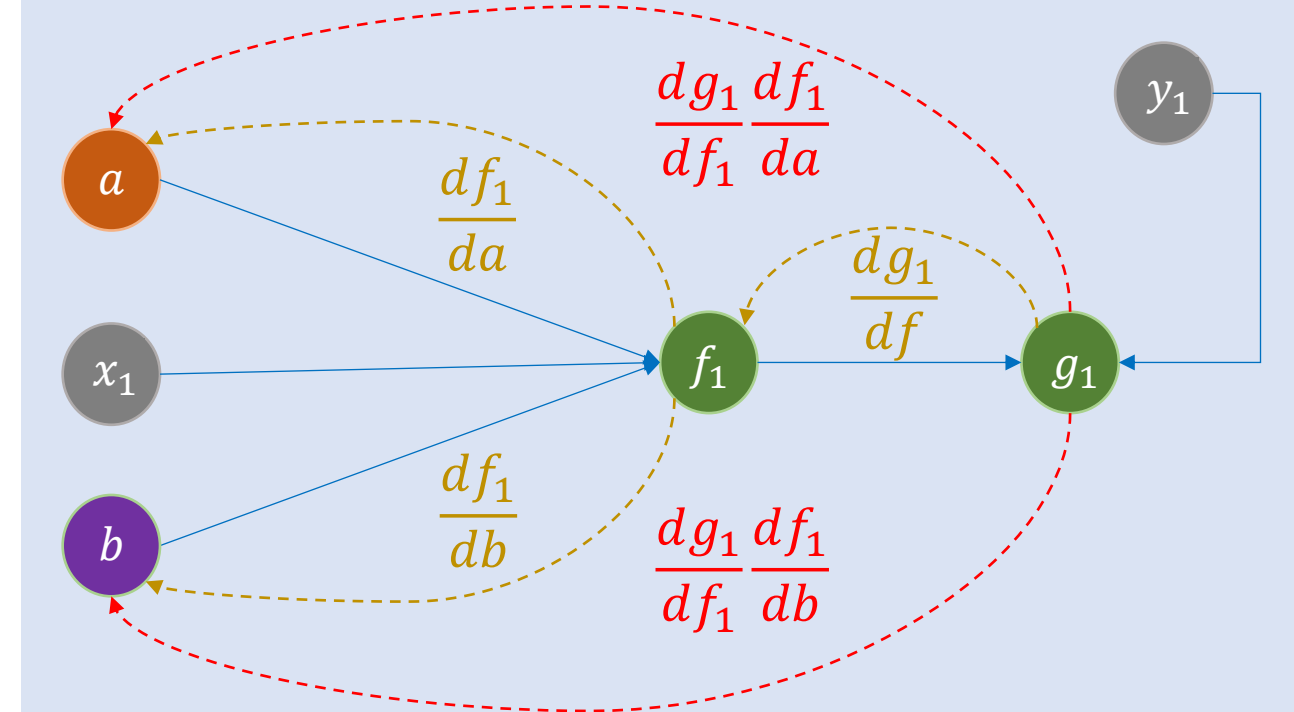
$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$

$$\frac{dg_1}{df_1} \frac{df_1}{da}$$

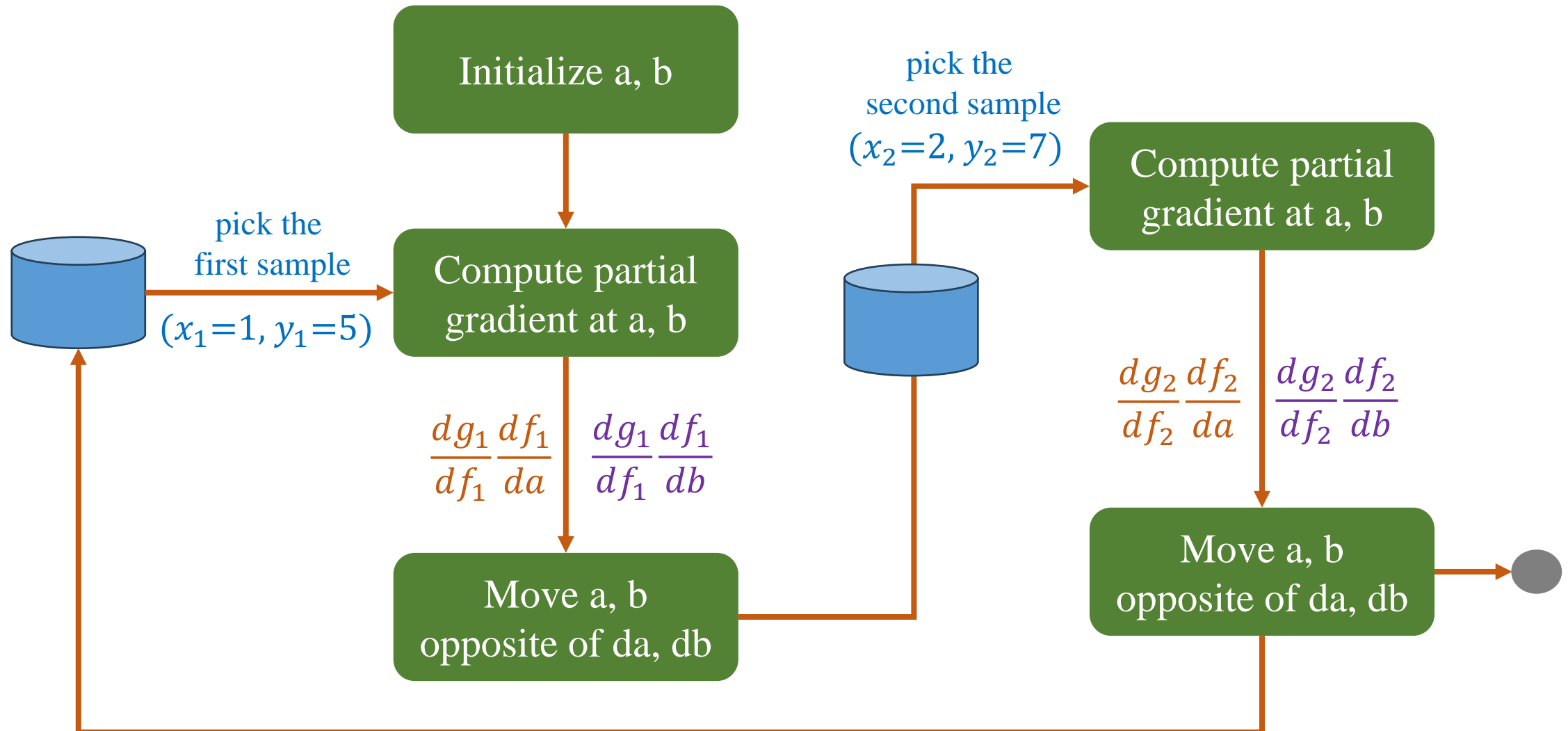
$$\frac{dg_2}{df_2} \frac{df_2}{da}$$

$$\frac{dg_1}{df_1} \frac{df_1}{db}$$

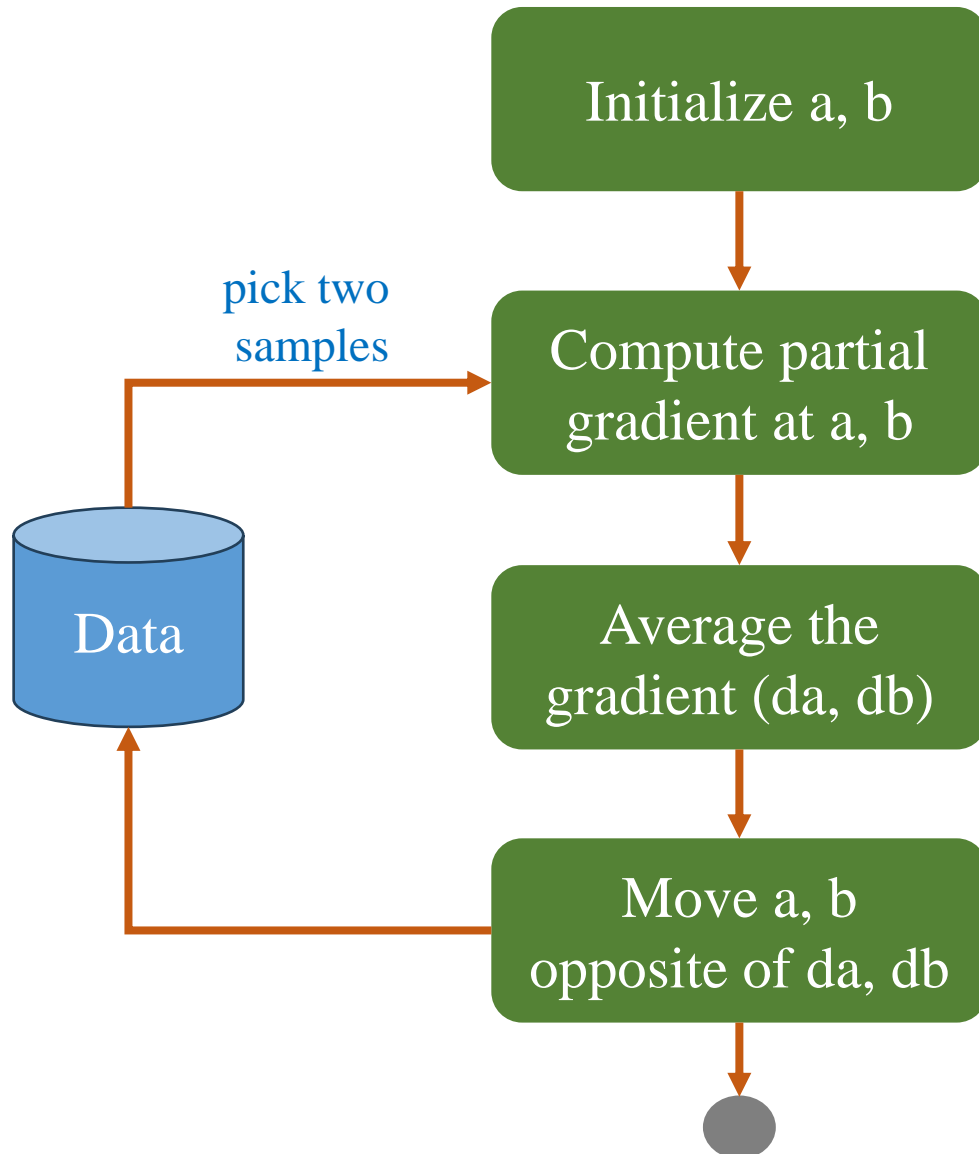
$$\frac{dg_2}{df_2} \frac{df_2}{db}$$



❖ Discussion: Approach 1



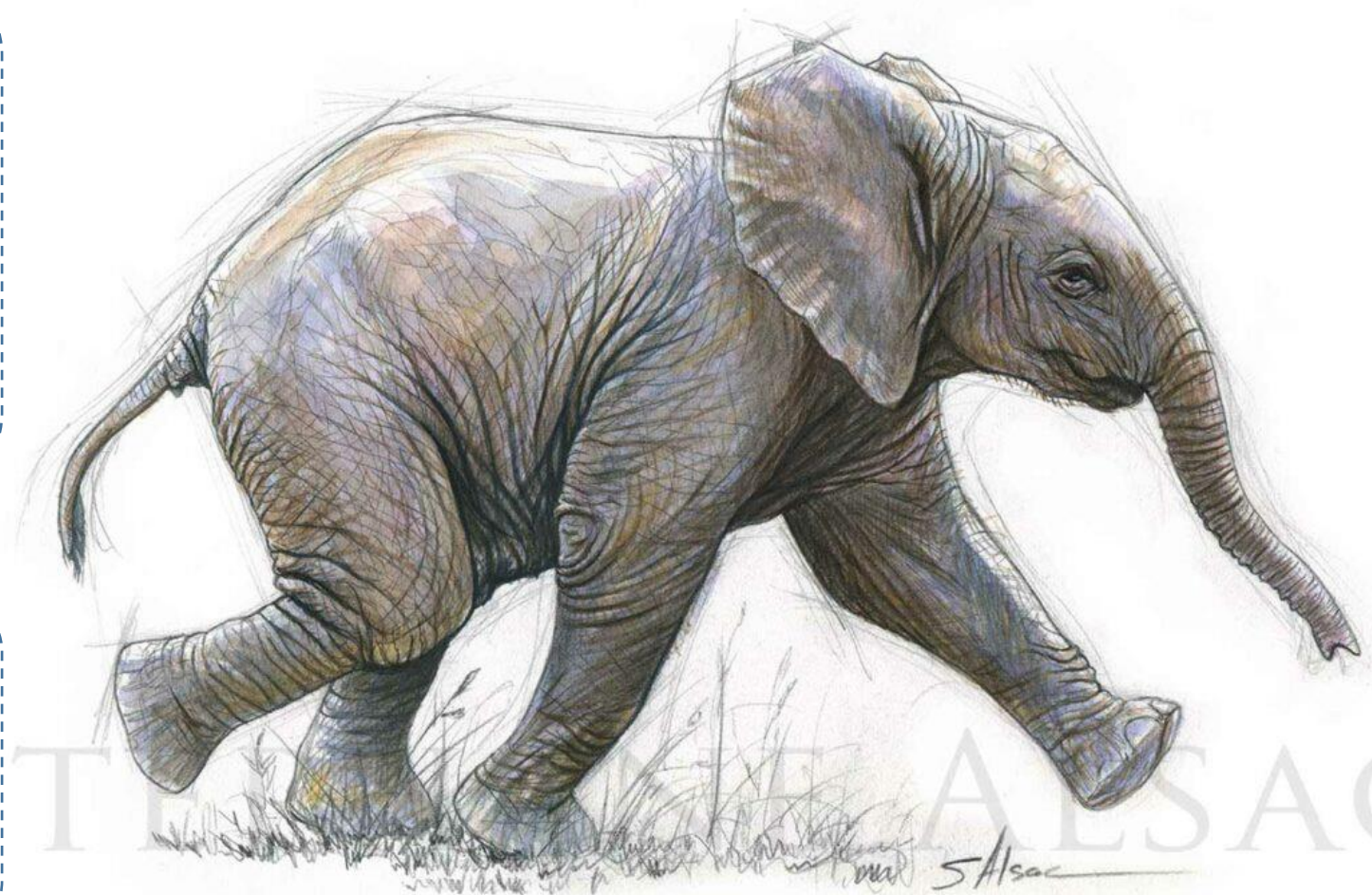
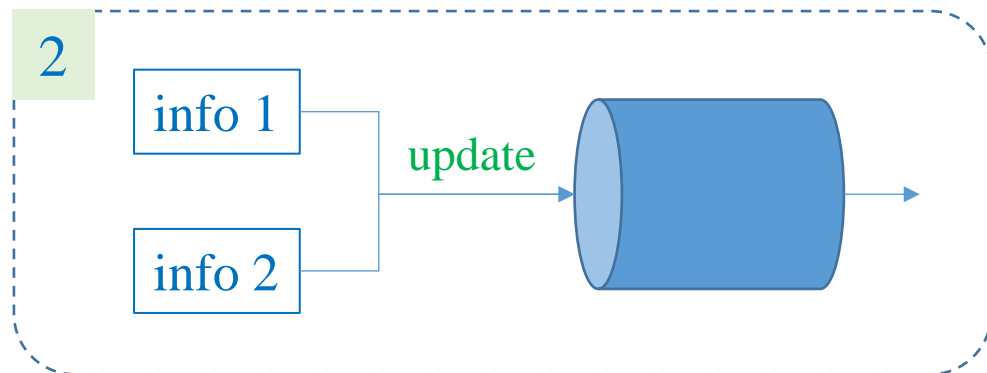
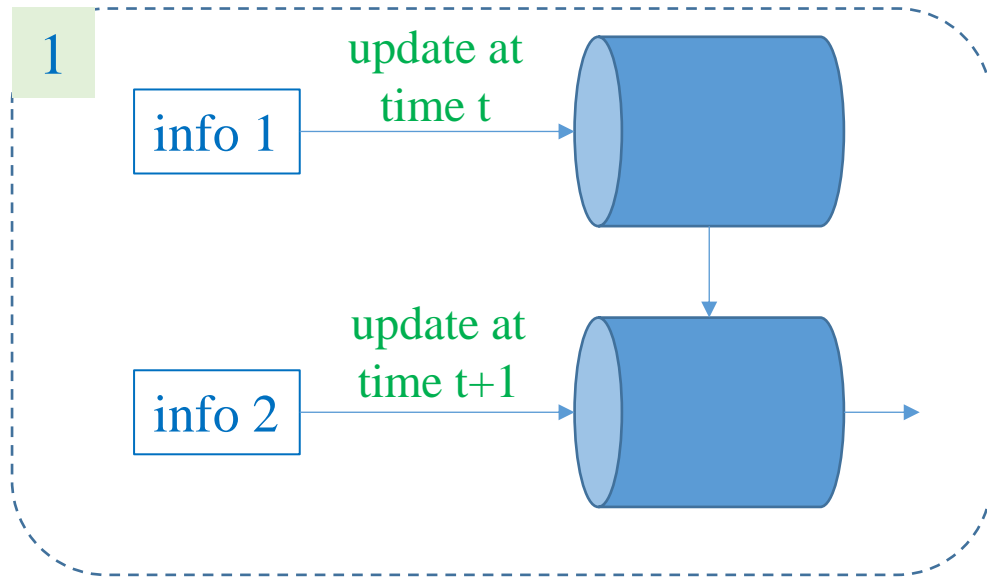
❖ Discussion: Approach 2



$$\frac{1}{2} \sum_i \frac{dg_i}{da} = \frac{1}{2} \left(\frac{dg_1}{df_1} \frac{df_1}{da} + \frac{dg_2}{df_2} \frac{df_2}{da} \right)$$

$$\frac{1}{2} \sum_i \frac{dg_i}{db} = \frac{1}{2} \left(\frac{dg_1}{df_1} \frac{df_1}{db} + \frac{dg_2}{df_2} \frac{df_2}{db} \right)$$

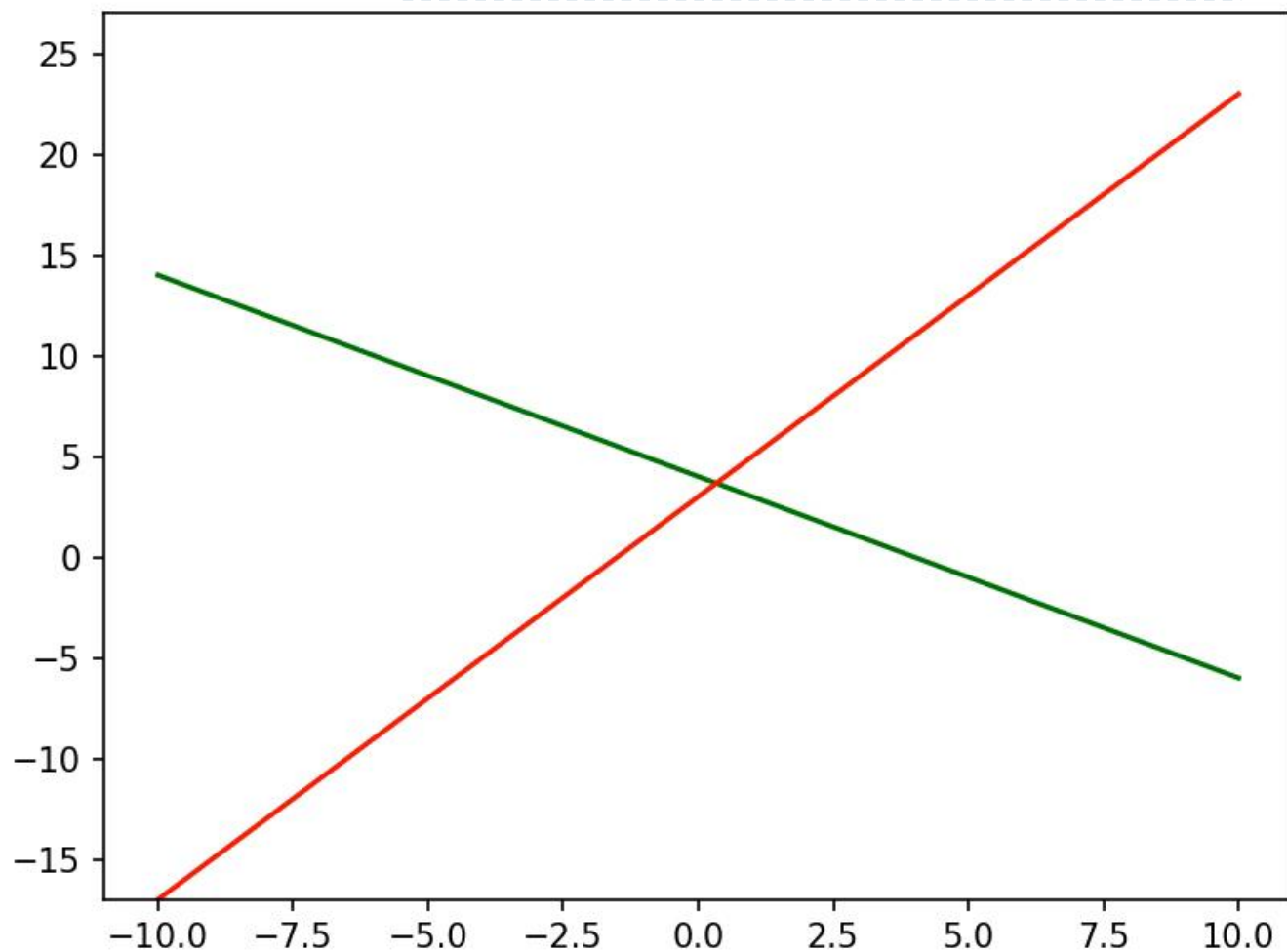
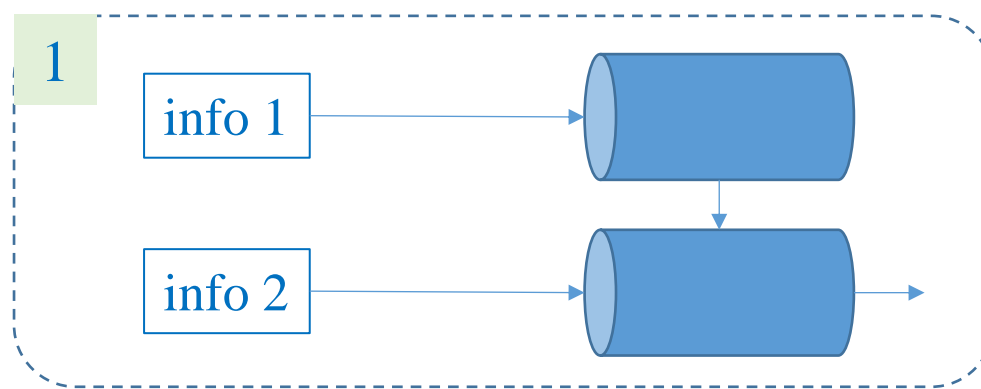
❖ How to use gradient information



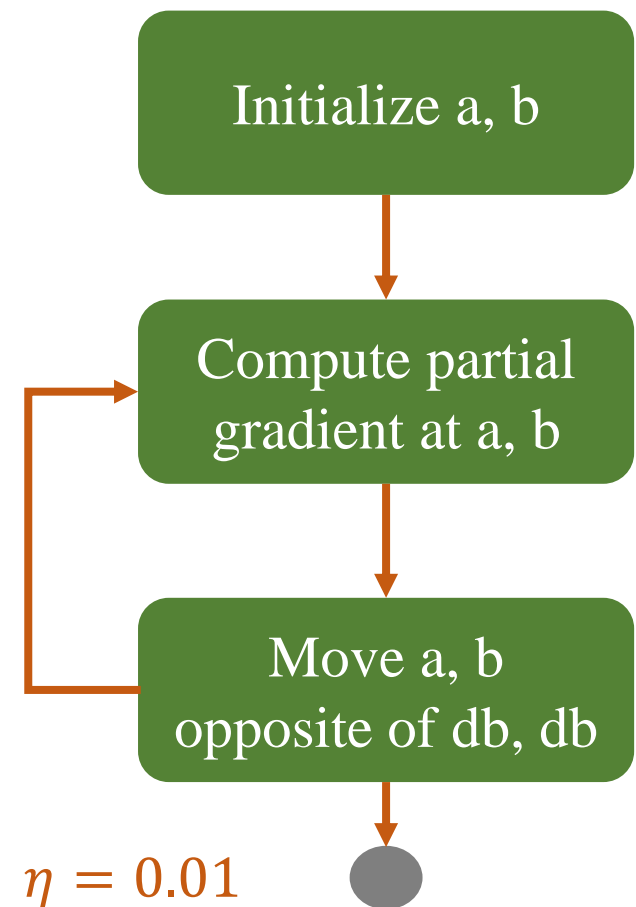

```
1 import random
2
3 # predict function
4 def predict_func(x, a, b):
5     return a*x + b
6
7 # parameters
8 num_steps = 100
9 lr = 0.01
10
11 # given data
12 data = [[-2, -1],
13         [5, 13]]
14
15 # 1. set a, b randomly
16 a = random.random()*10.0 - 5.0
17 b = random.random()*10.0 - 5.0
```

```
1 for i in range(num_steps):
2     for sample in data:
3         x_value, y_value = sample
4
5         # compute predicted_y
6         predicted_y = predict_func(x_value, a, b)
7
8         # compute g
9         g_value = (predicted_y - y_value)**2
10
11         # compute partial gradients for a and b
12         dg_da = 2*x_value*(predicted_y - y_value)
13         dg_db = 2*(predicted_y - y_value)
14
15         # update
16         a = a - lr*dg_da
17         b = b - lr*dg_db
```


Summary



$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$
$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$

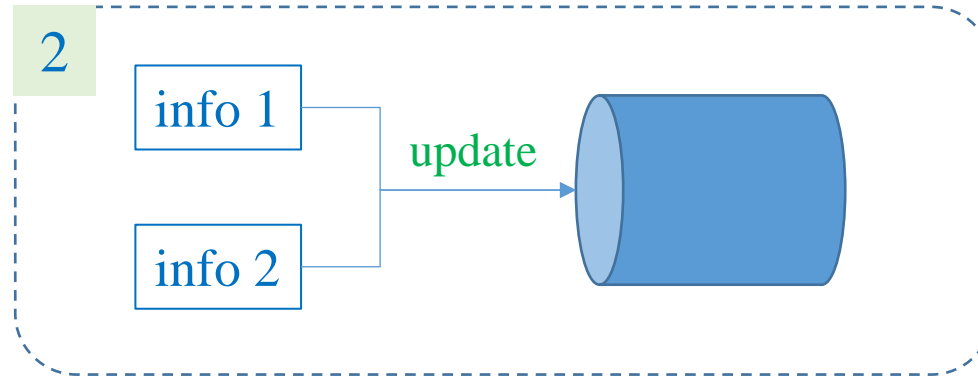


```
1 import random
2
3 # predict function
4 def predicted_func(x, a, b):
5     return a*x + b
6
7 # parameters
8 num_steps = 100
9 lr = 0.01
10
11 # given data
12 x1 = -2
13 y1 = -1
14
15 x2 = 5
16 y2 = 13
17
18 # 1. set a, b randomly
19 a = random.random()*10.0 - 5.0
20 b = random.random()*10.0 - 5.0
```

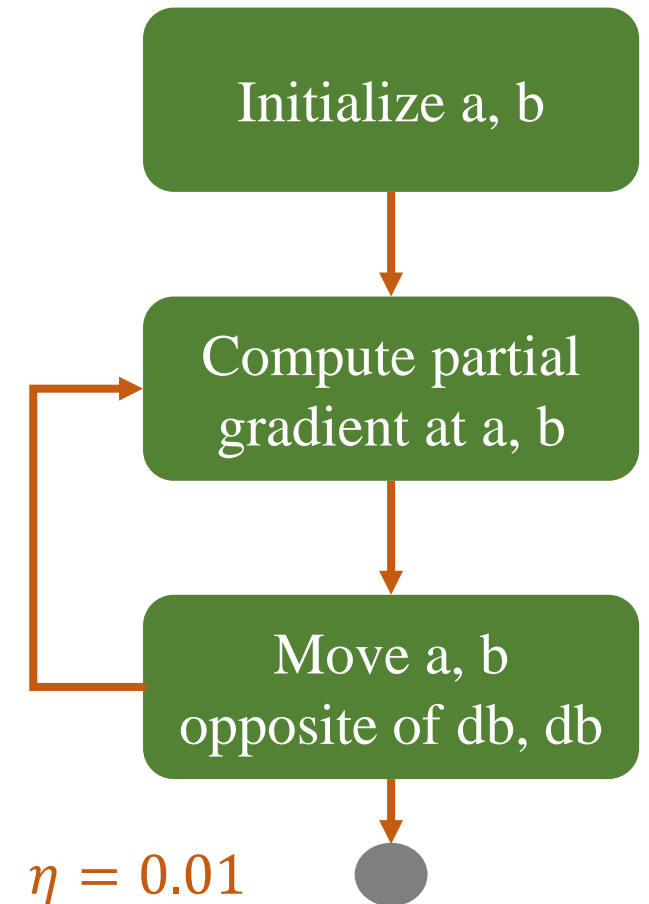
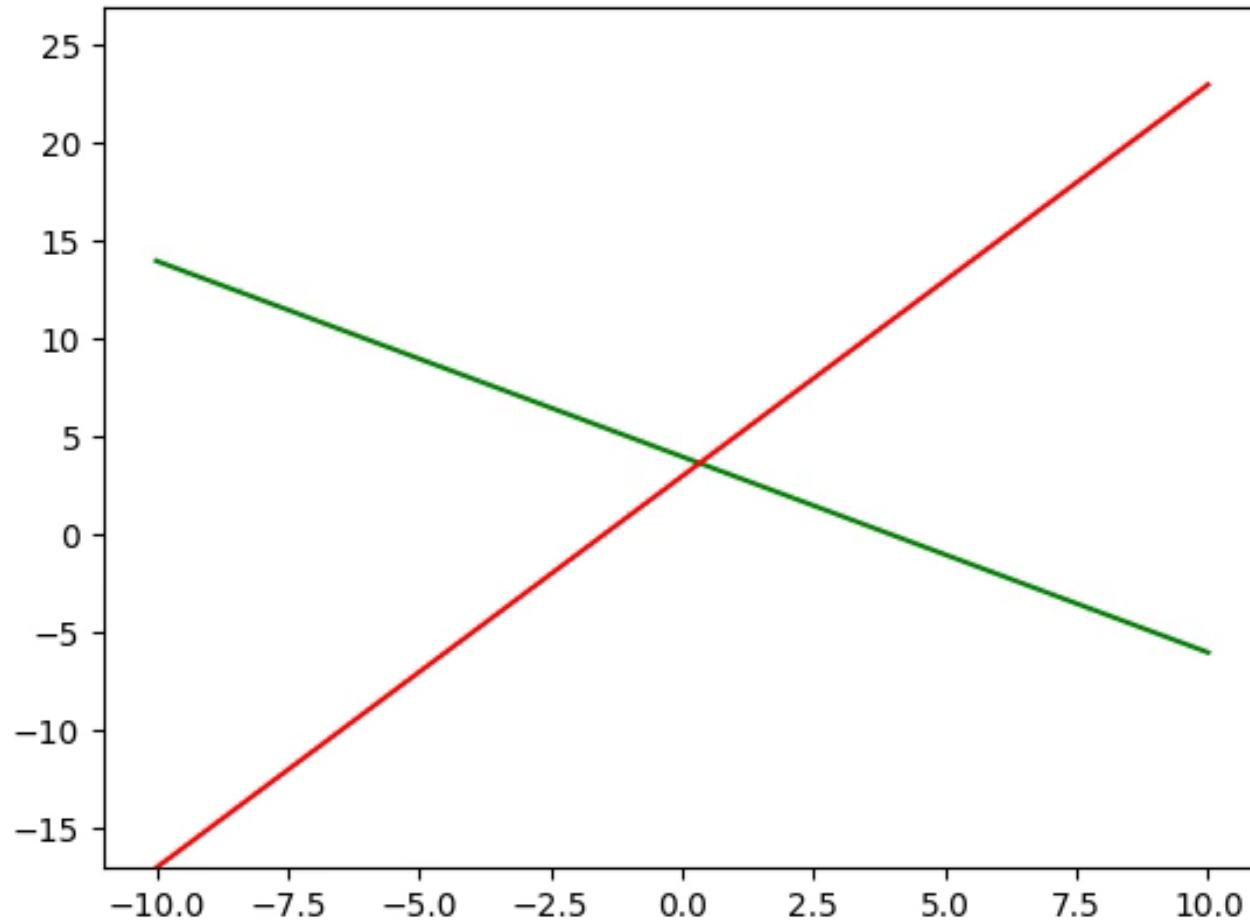
```
1 for i in range(num_steps):
2     # 2. compute predicted_y1 and predicted_y2
3     predicted_y1 = predicted_func(x1, a, b)
4     predicted_y2 = predicted_func(x2, a, b)
5
6     # 3. compute g
7     g_value_1 = (predicted_y1 - y1)**2
8     g_value_2 = (predicted_y2 - y2)**2
9
10    # logging
11    # ...
12
13    # 4. compute partial gradients for a and b
14    dg_da = 2*x1*(predicted_y1 - y1) + 2*x2*(predicted_y2 - y2)
15    dg_db = 2*(predicted_y1 - y1) + 2*(predicted_y2 - y2)
16
17    # 5. update
18    a = a - lr*dg_da/2
19    b = b - lr*dg_db/2
```

Summary

2



$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$
$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$



Summary

2

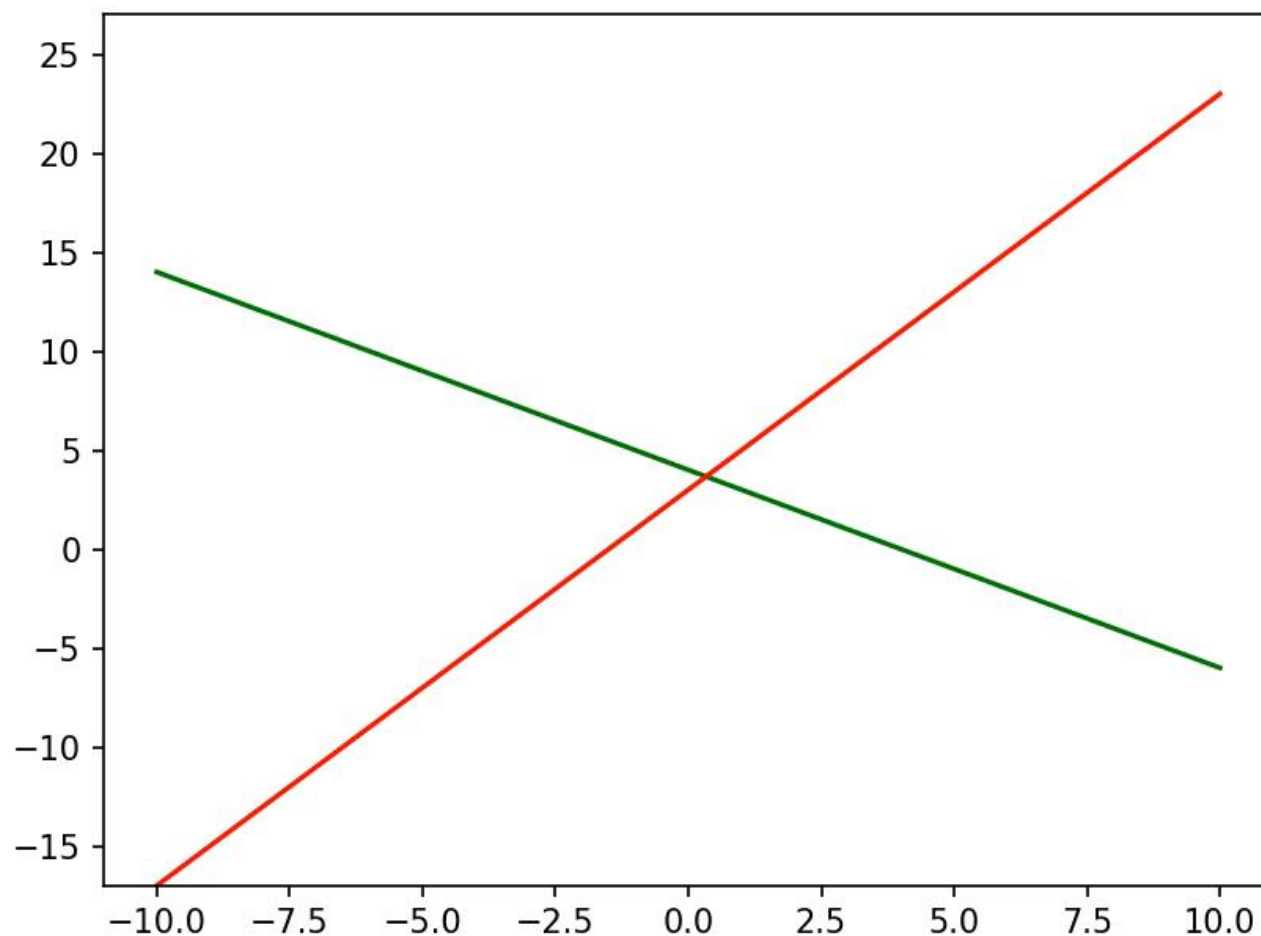
info 1

info 2

update



$$\frac{dg}{da} = \frac{dg}{df} \frac{df}{da} = 2x(f - y)$$
$$\frac{dg}{db} = \frac{dg}{df} \frac{df}{db} = 2(f - y)$$



Initialize a, b

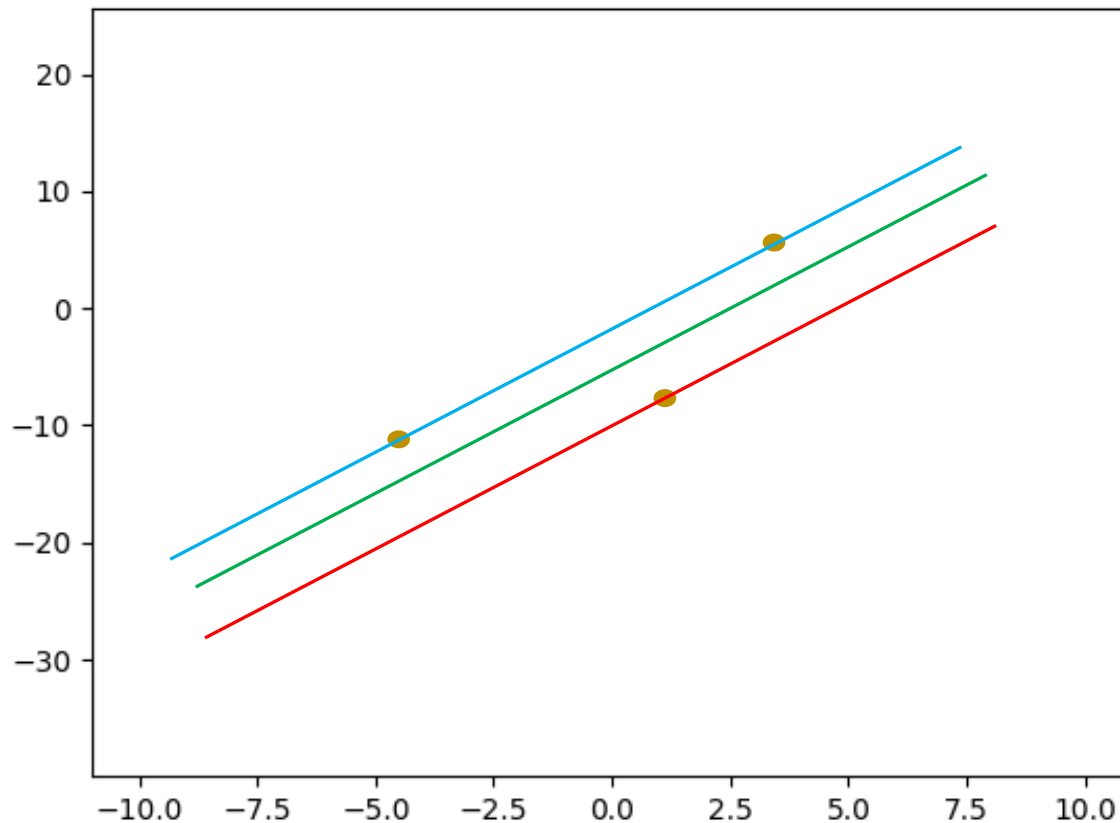
Compute partial
gradient at a, b

Move a, b
opposite of db, db

$\eta = 0.001$



❖ What about this given samples?



Which line is the best representation of these three points?



Line 1: go through two points

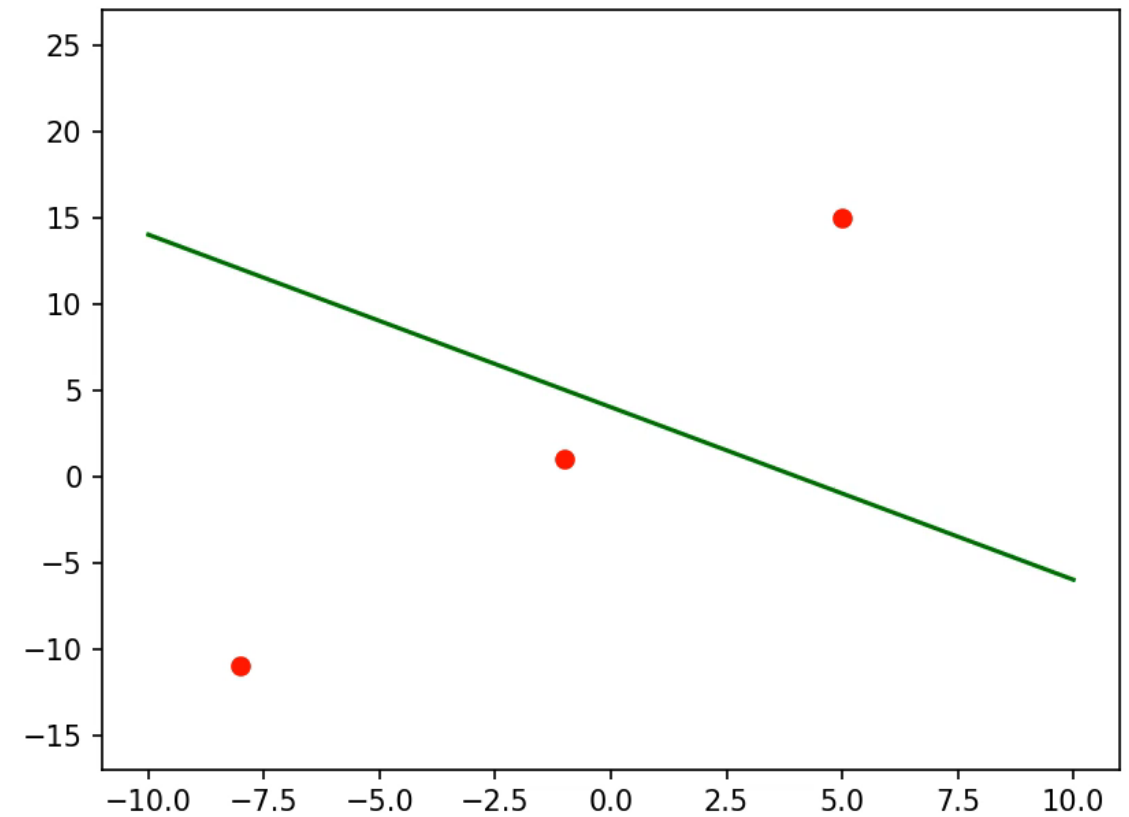
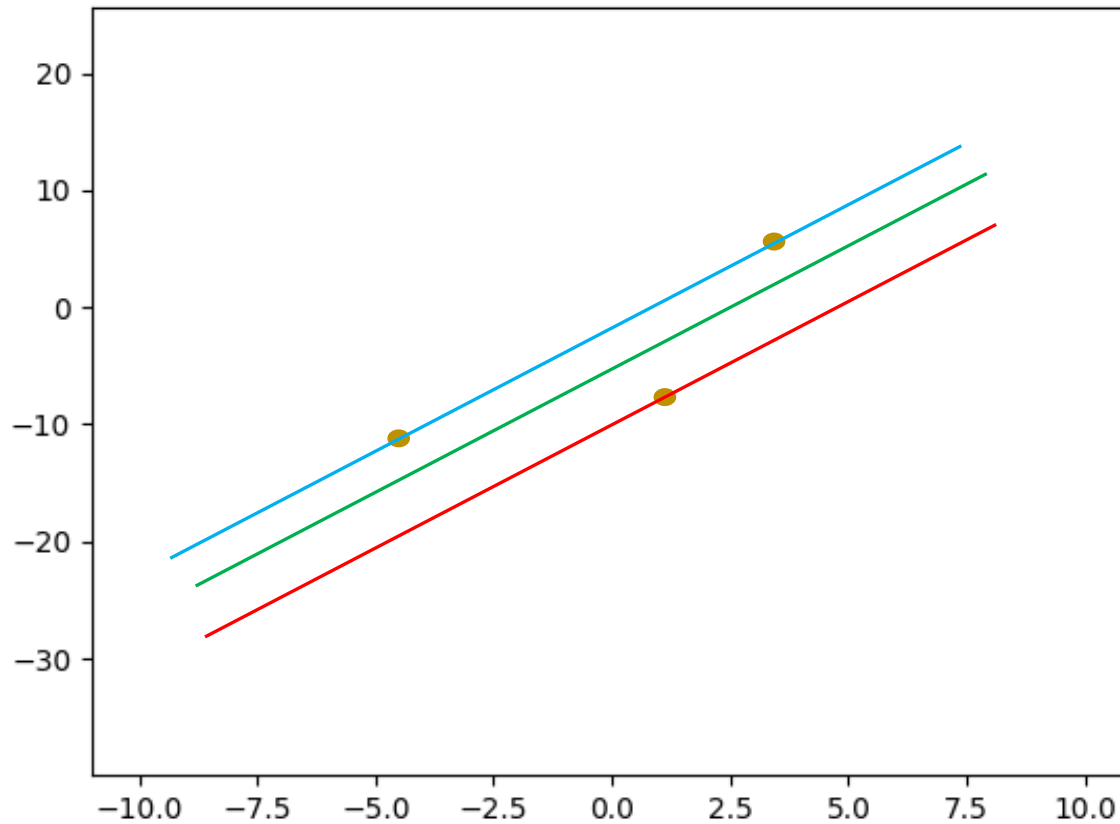


Line 2: smallest summation of distances



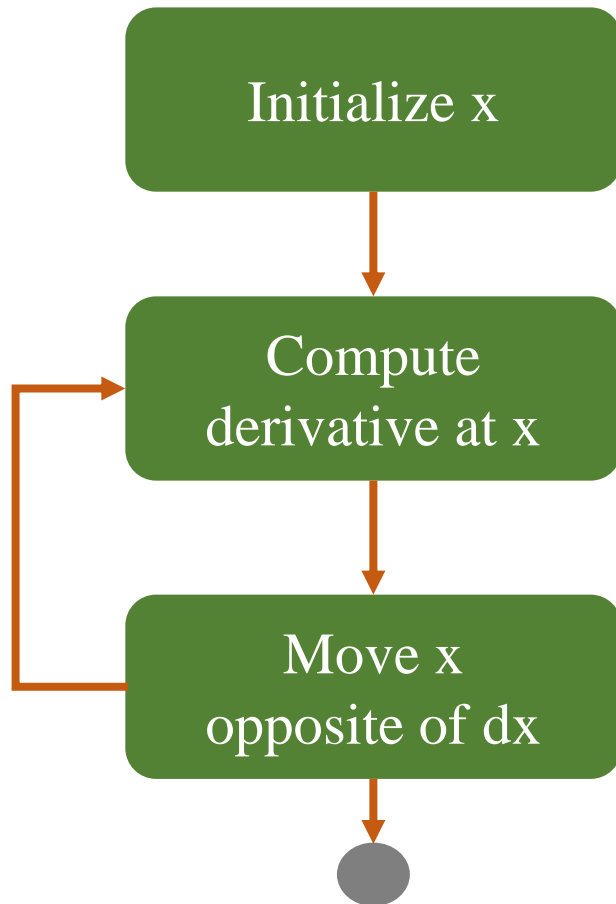
Line 3: go through one point

❖ What about the given samples?

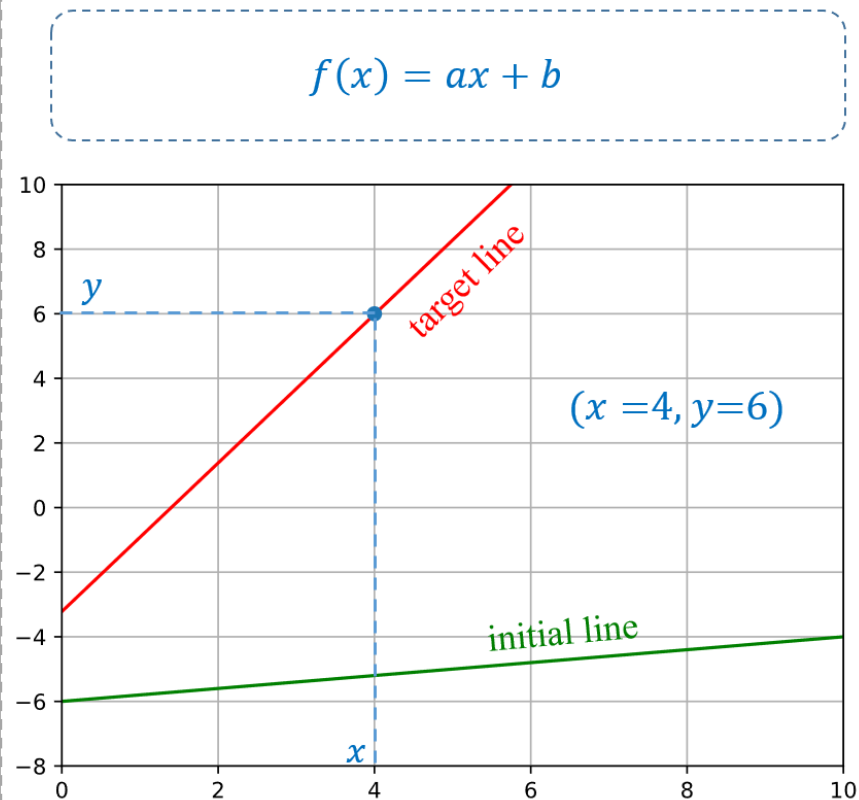


Objectives

Optimization



Problem Solving



Toward Linear Regression

