# Steps into Linear Regression

Quang–Vinh Dinh
PhD in Computer Science

*Year 2023*

# **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

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

House price data

$$price = w * area + b$$

if area=6.0, price=7.28

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

House price data

$$price = w_1 * area + b_1$$
$$price = w_2 * area + b_2$$
$$price = w_3 * area + b_3$$

if area=6.0, price=?

# Outline

- ➢ **Review on Optimization**
- ➢ **Partial Gradient and Chain Rules**
- ➢ **Finding a Line**
- ➢ **Linear Regression**

# Optimization

❖ **Gradient descent**



$$\frac{d}{dx}J(x) = \lim_{\Delta x \to 0} \frac{J(x + \Delta x) - J(x)}{\Delta x}$$
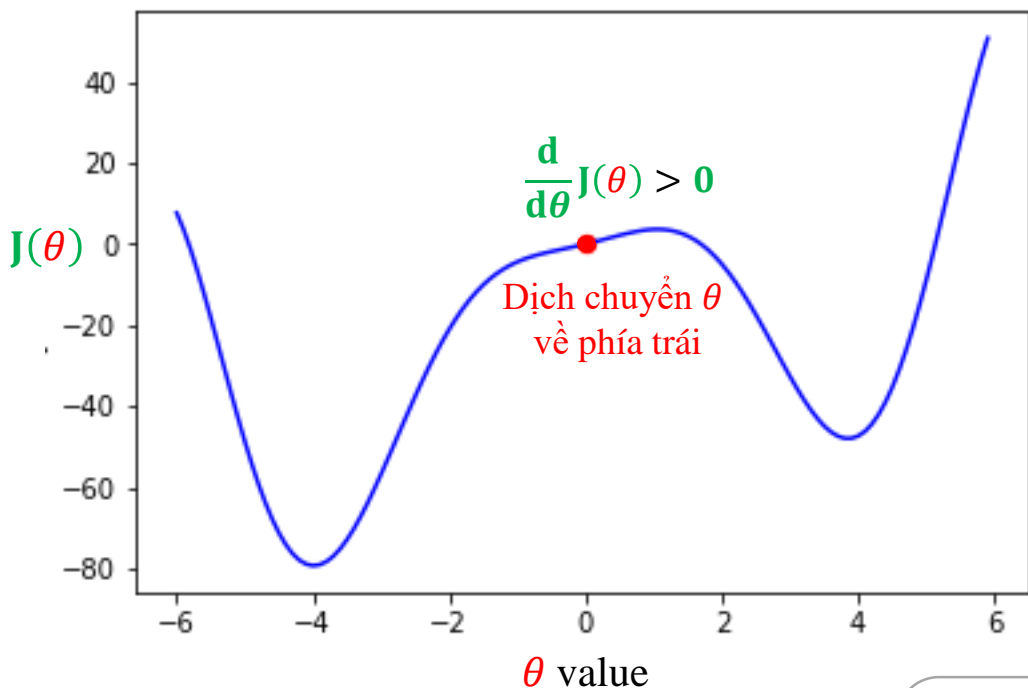
$$\frac{d}{dx}J(x_2) < 0 \qquad \frac{d}{dx}J(x_1) > 0$$

$$\mathbf{x} = \mathbf{x} - \eta \frac{d}{dx}\mathbf{J(x)}$$

Derivate at x

learning rate

3

# Optimization

❖ **Gradient descent**

**Khởi tạo giá trị $\theta$**



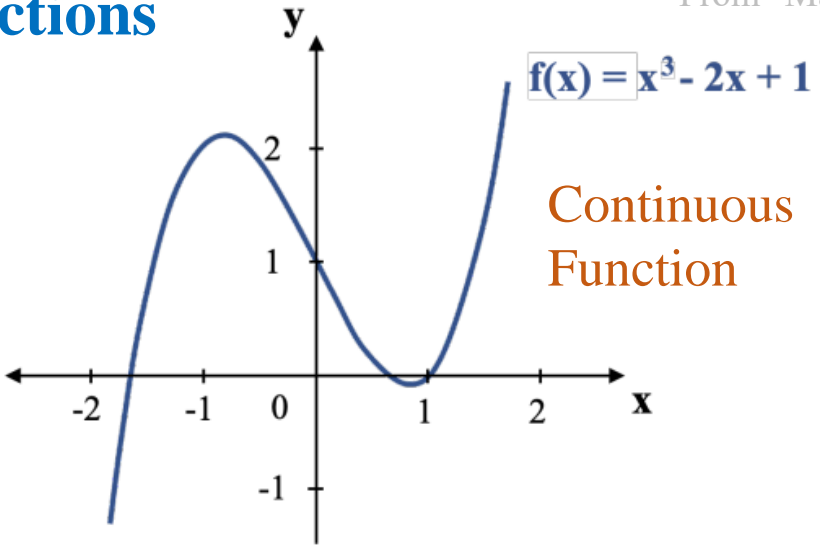$\frac{d}{d\theta}J(\theta) > 0$

Dịch chuyển $\theta$
về phía trái

**Di chuyển $\theta$ ngược hướng đạo hàm**



$\frac{d}{d\theta}J(\theta) > 0$

Dịch chuyển $\theta$
về phía trái

$\frac{d}{d\theta}J(\theta) > 0$

Dịch chuyển $\theta$
về phía trái

**Cứ tiếp tục di chuyển ngược hướng đạo hàm**



$\frac{d}{d\theta}J(\theta) > 0$

Dịch chuyển $\theta$ về phía trái

$\frac{d}{d\theta}J(\theta) < 0$

Dịch chuyển $\theta$
về phía phải

4

# Optimization Algorithms

## Loss functions

From "Machine Learning Simplified"

$$f(x) = \begin{cases} 3x^2 - 2, & x < 1 \\ 2x - 1, & x > 1 \end{cases}$$

f(x) = 2x - 1

f(x) = x³ - 2x + 1

f(x) = x³ - 2x + 1

Continuous Function
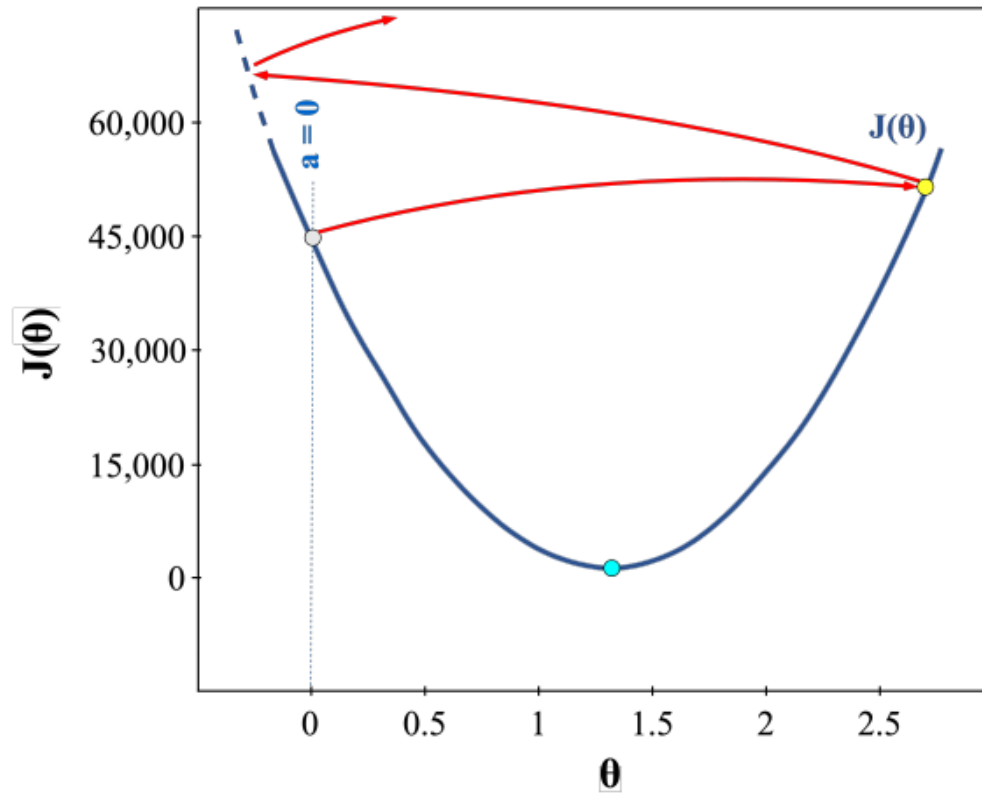
f(x) = x^{1/3}

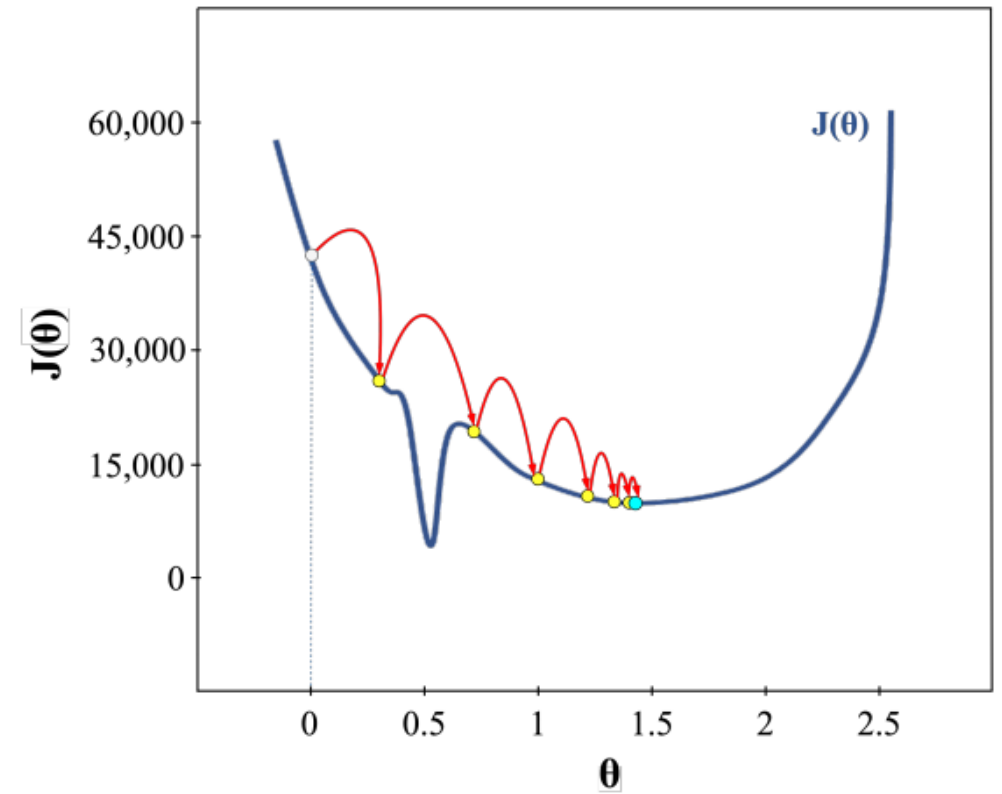Continuous non-differentiable functions

f(x) = 1/x

Discontinuous Functions

5

# Optimization Algorithms

## Learning rate



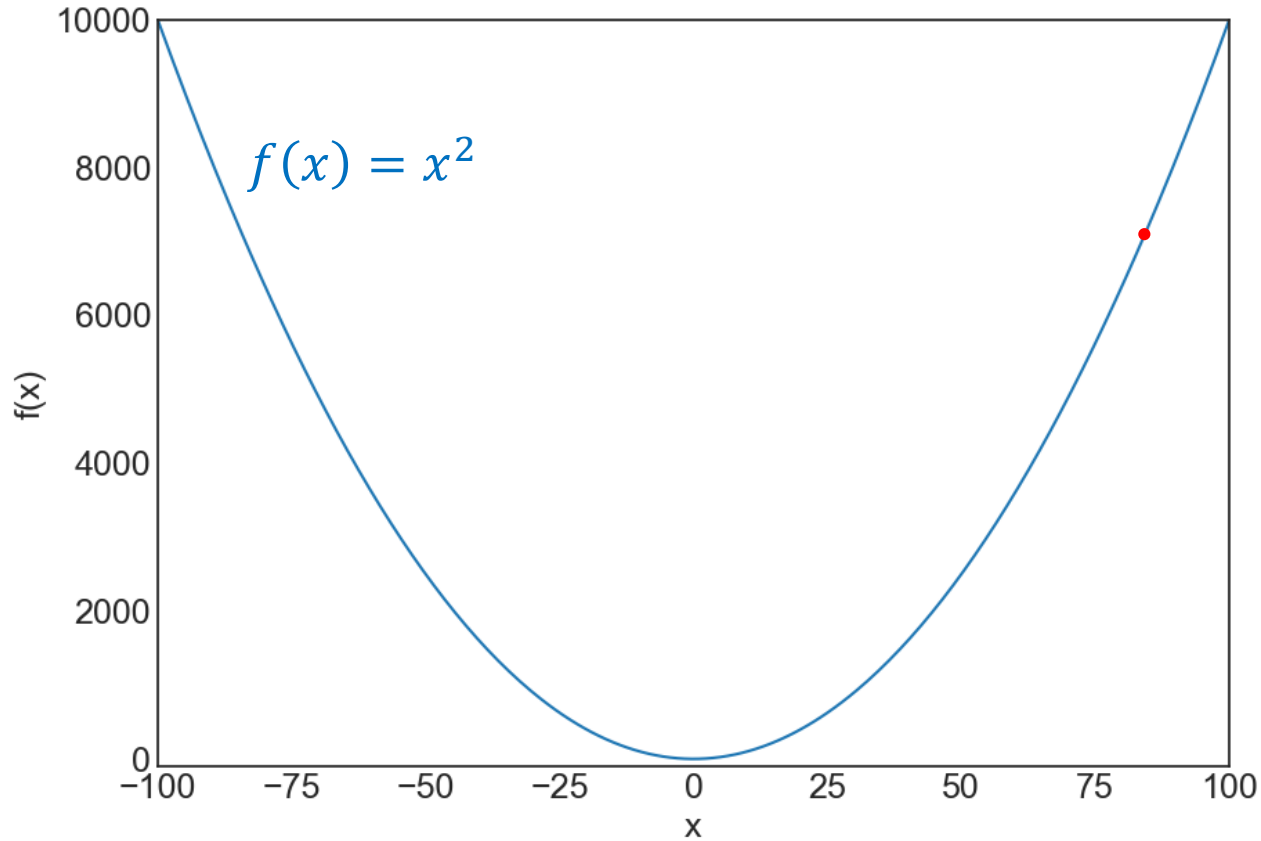(a) Gradient descent missing global minimum on a convex cost function due to a very large learning rate.

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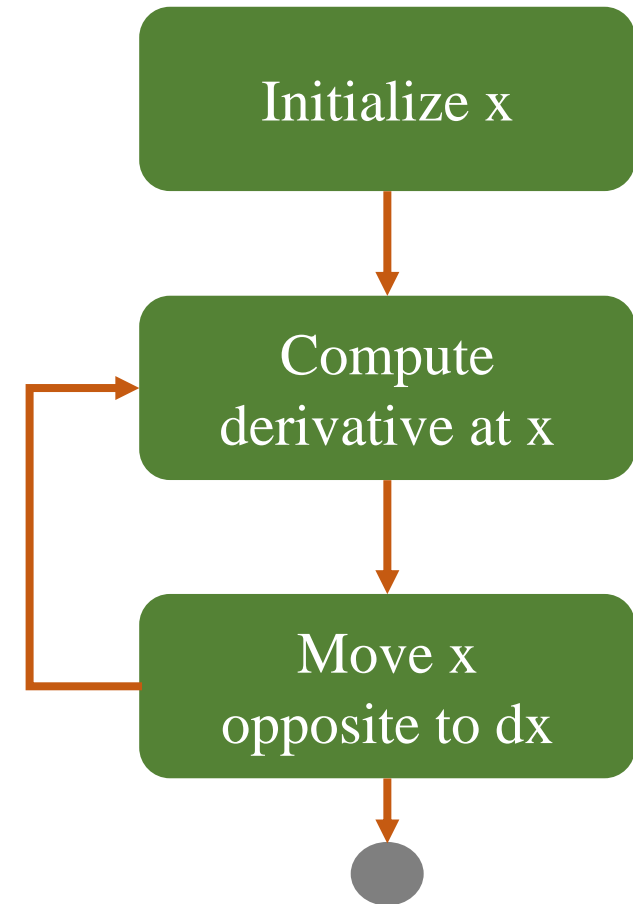
From "Machine Learning Simplified"

6

# Gradient-based Optimization

❖ **Square function**



$f(x) = x^2$

$$-100 \leq x \leq 100$$
$$x \in \mathbb{N}$$

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

Initialize x

Compute derivative at x

Move x opposite to dx

**7**

# Optimization

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

❖ **Square function**



$f(x) = x^2$

$-100 \leq x \leq 100$

$x \in \mathbb{N}$

$x_0 = 70.0 \qquad \eta = 0.1$

$f'(x_0) = 140.0$

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

$f'(x_1) = 112.0$

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

$f'(x_2) = 89.6$

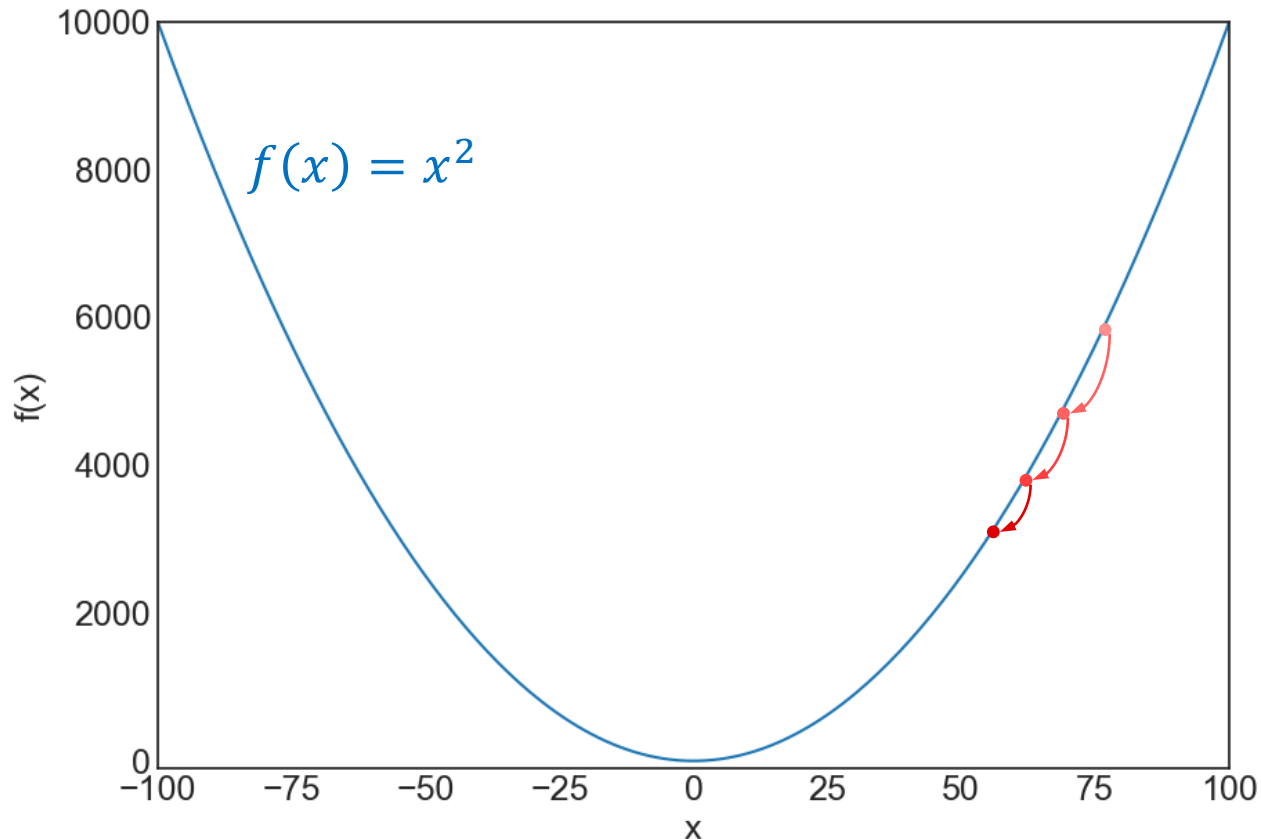$x_3 = x_2 - \eta f'(x_2) = 35.84$

$f'(x_3) = 71.68$

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

# Optimization

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

❖ **Square function**



$f(x) = x^2$

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

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

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

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

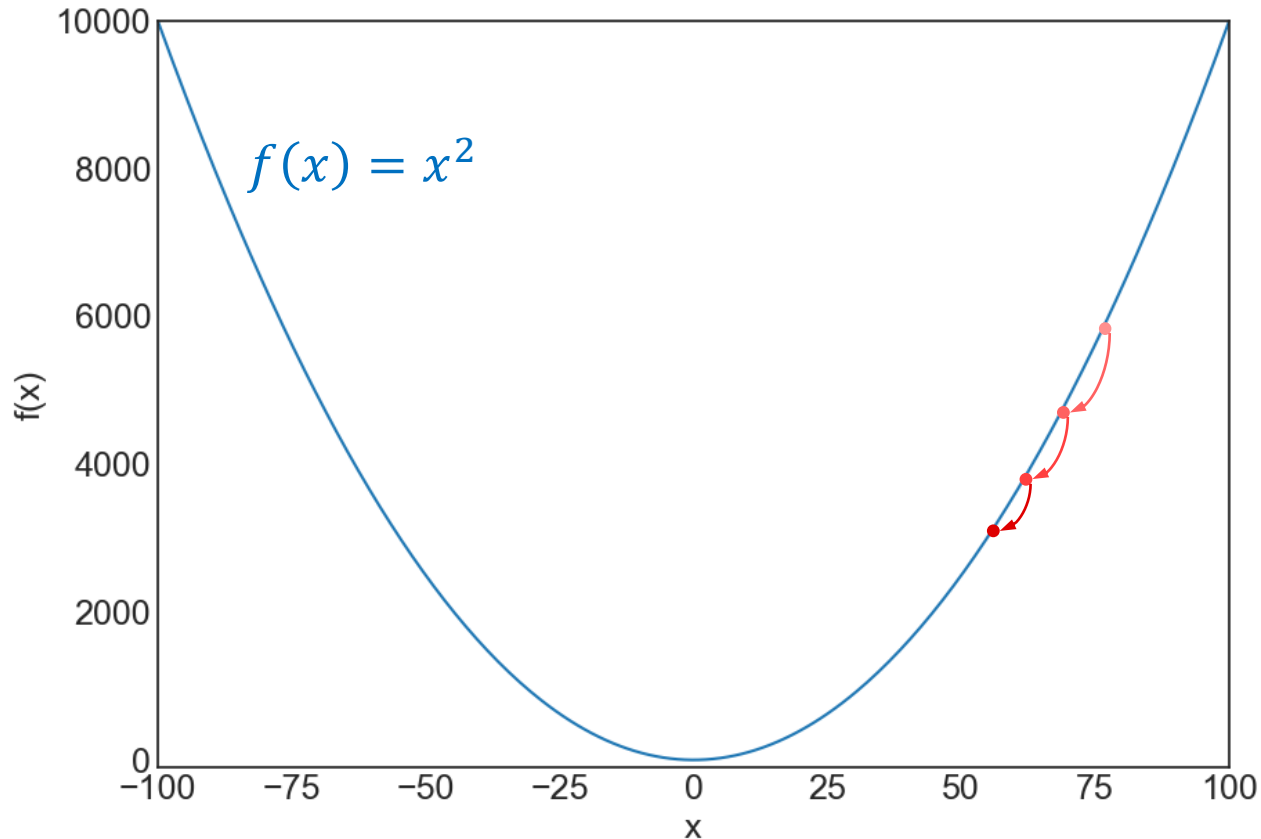$x_{12} = x_{11} - \eta f'(x_{11}) = 3.84$

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

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

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

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

Keep doing

9

# **Optimization**

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

❖ **Square function**



$f(x) = x^2$

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

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

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

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

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

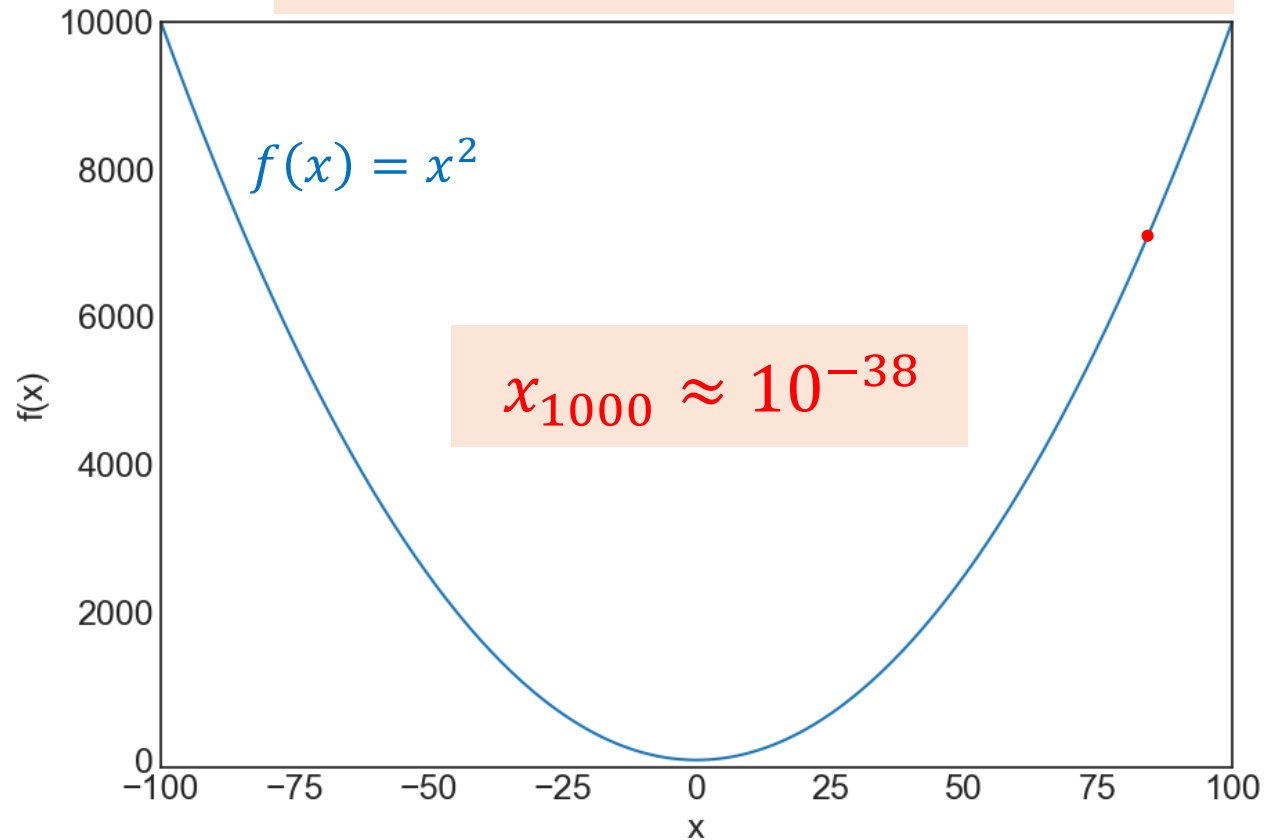$f'(x_{32}) = 0.88$

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

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

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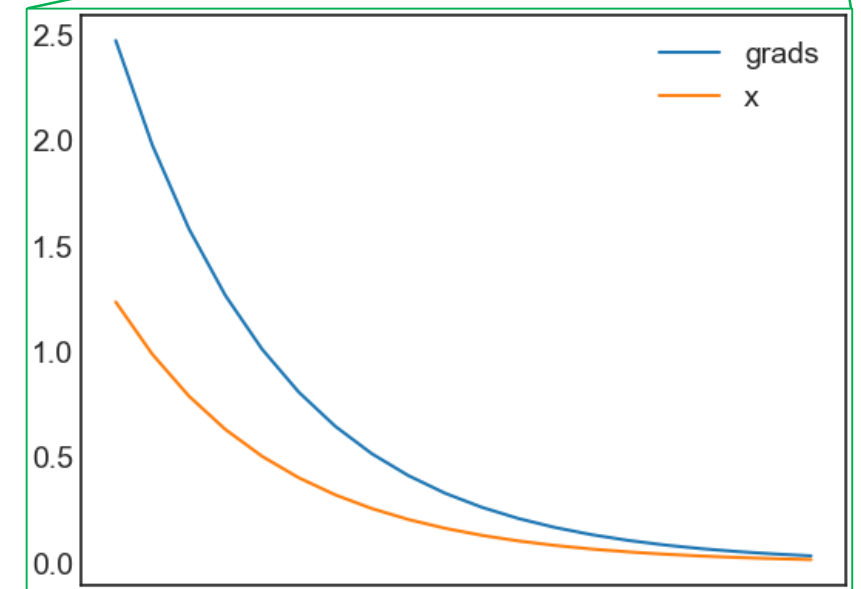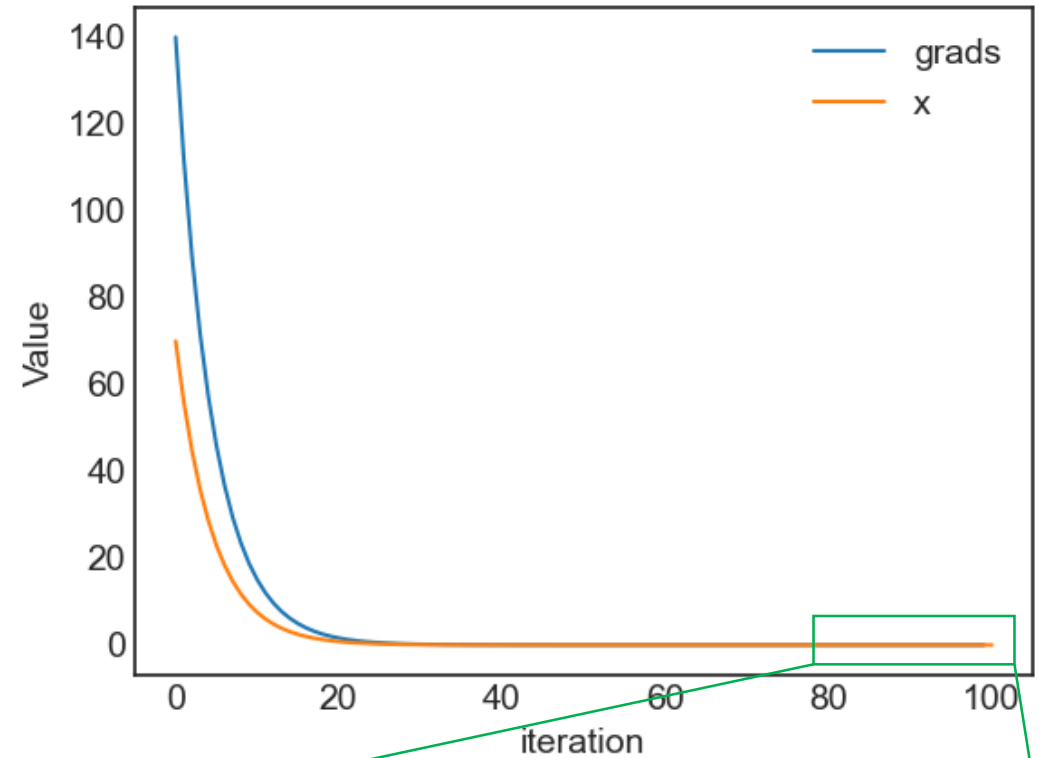
Keep doing

# Optimization

❖ **Square function**

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

$f(x) = x^2$

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

Optimized successfully!

# Optimization

❖ **Square function**

$x_0 = 99.0$

$\eta = 0.1$

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



SGD

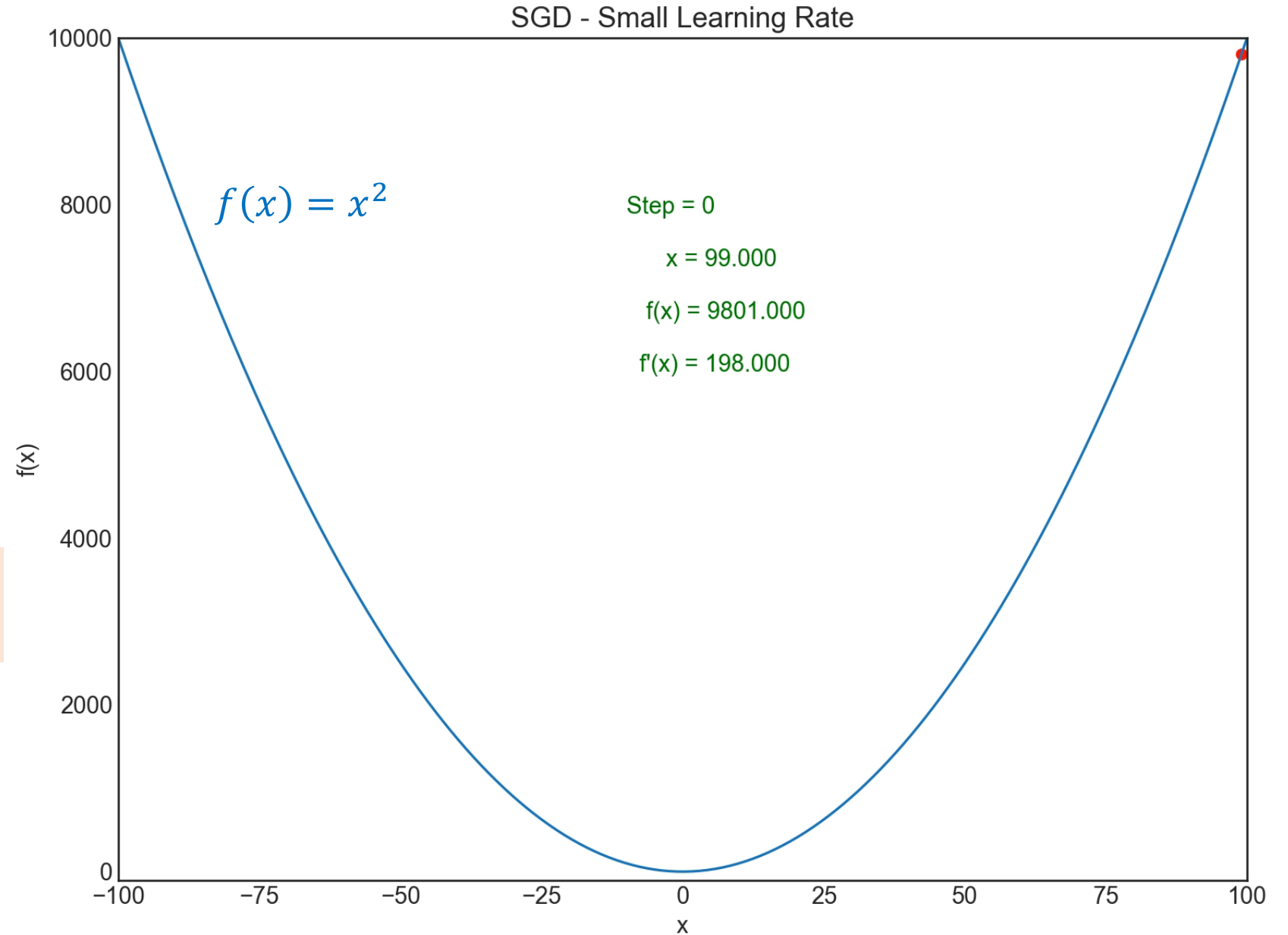$f(x) = x^2$

Step = 1

x = 79.200

f(x) = 6272.640

f'(x) = 158.400

# Optimization

❖ **Square function**

Discussion

$$x_0 = 99.0$$

$$\eta = 0.001$$

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

SGD - Small Learning Rate

$f(x) = x^2$

Step = 0

x = 99.000

f(x) = 9801.000

f'(x) = 198.000

# Optimization

❖ **Square function**

Discussion

$$x_0 = 99.0$$

$$\eta = 0.8$$

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



SGD - Large Learning Rate

$f(x) = x^2$

Step = 0

x = 99.000

f(x) = 9801.000

f'(x) = 198.000

# Optimization

❖ **Square function**

<span style="color:red">Discussion</span>

$$x_0 = 99.0$$

$$\eta = 1.1$$

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



SGD - Too Large Learning Rate

$f(x) = x^2$

Step = 0

x = 5.000

f(x) = 25.000

f'(x) = 10.000

15

# Optimization

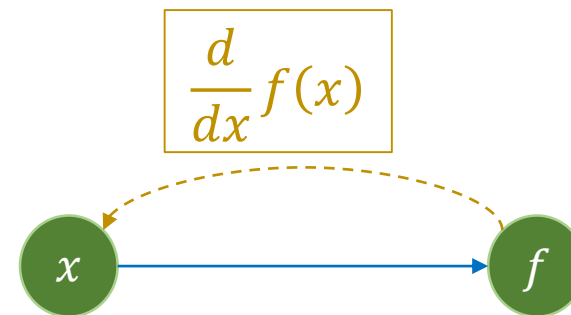$$x_t = x_{t-1} - \eta f'(x_{t-1})$$

❖ **Square function**



$f(x) = x^2$

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

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

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

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

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

$$\frac{d}{dx}f(x)$$



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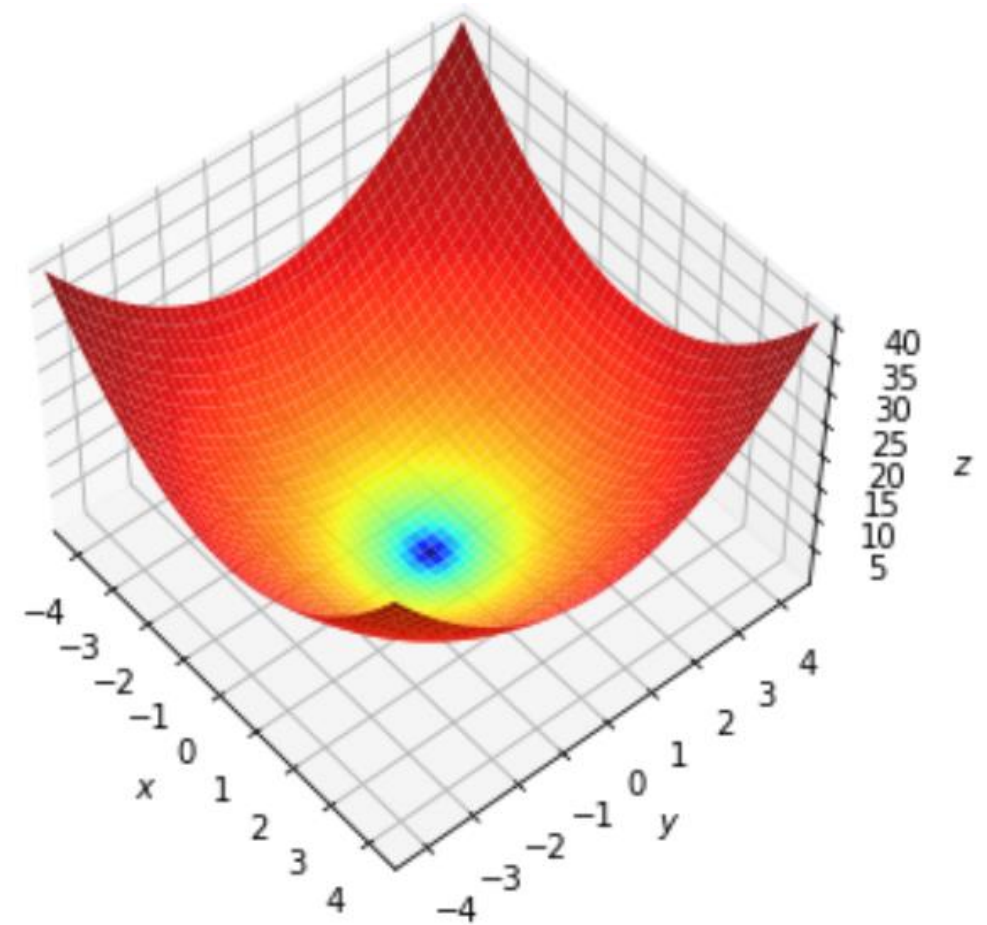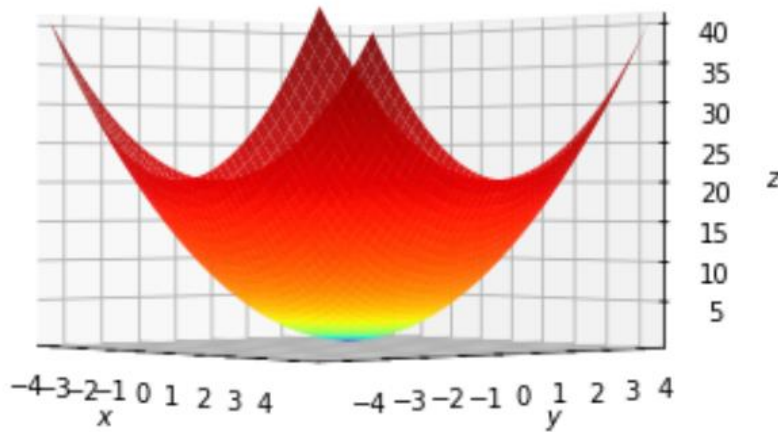- After an update, $f(x_{new}) \leq f(x_{old})$

16

# Optimization

❖ **Optimization: 2D function**

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

$$-100 \leq x, y \leq 100$$
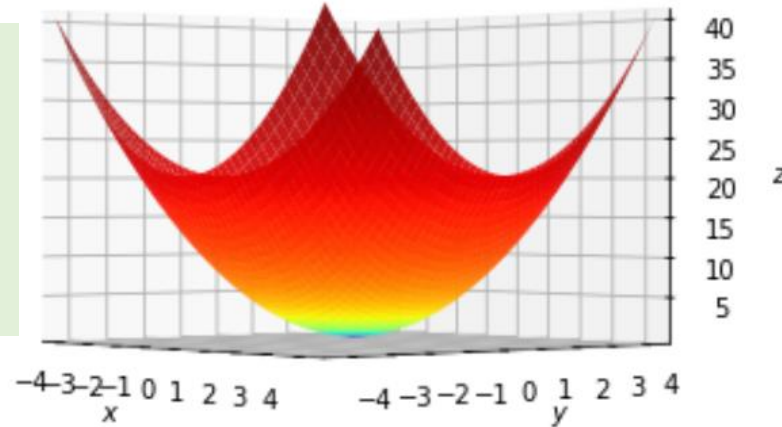
$$x, y \in \mathbb{N}$$

# Optimization

## ❖ Optimization: 2D function

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

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

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



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

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

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

$\dfrac{\partial f(x_0, y_0)}{\partial x} = 12 \qquad \dfrac{\partial f(x_0, y_0)}{\partial y} = 18$

$x_1 = 4.8 \qquad\qquad y_1 = 7.2$

$\dfrac{\partial f(x_1, y_1)}{\partial x} = 9.6 \qquad \dfrac{\partial f(x_1, y_1)}{\partial y} = 14.4$

$x_2 = 3.84 \qquad\qquad y_2 = 5.75$

$\dfrac{\partial f(x_2, y_2)}{\partial x} = 7.68 \qquad \dfrac{\partial f(x_2, y_2)}{\partial y} = 11.51$

$x_3 = 3.07 \qquad\qquad y_3 = 4.608$

$\dfrac{\partial f(x_3, y_3)}{\partial x} = 6.14 \qquad \dfrac{\partial f(x_3, y_3)}{\partial y} = 9.21$

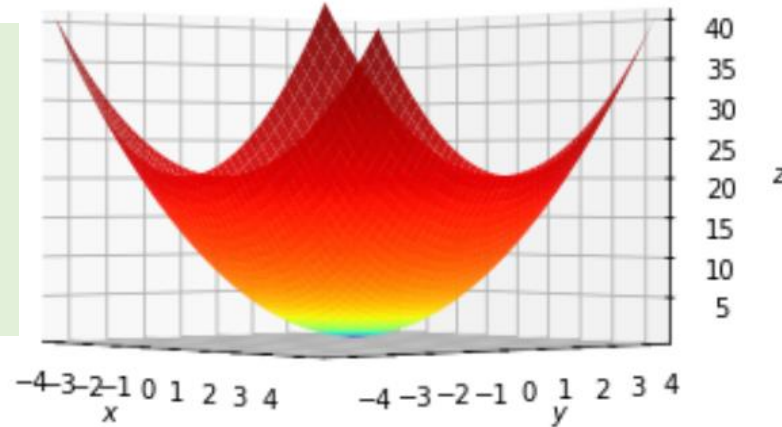$x_4 = 2.45 \qquad\qquad y_4 = 3.68$

# Optimization

❖ **Optimization: 2D function**

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

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

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



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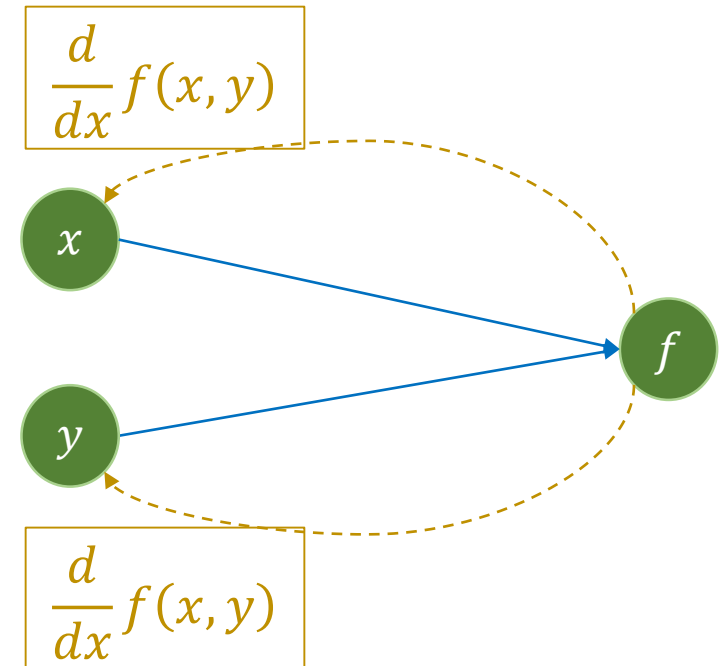
$$x_3 = 3.07 \qquad y_3 = 4.608$$

$$\frac{\partial f(x_3, y_3)}{\partial x} = 6.14 \qquad \frac{\partial f(x_3, y_3)}{\partial y} = 9.21$$

$$x_4 = 2.45 \qquad y_4 = 3.68$$

Summary:

- Given a function f(x, y), find optimal $(x_{opt}, y_{opt})$ so that $f(x_{opt}, y_{opt})$ is minimum

- After an update, $f(x_{new}, y_{new}) \leq f(x_{old}, y_{old})$

$$\frac{d}{dx} f(x, y)$$



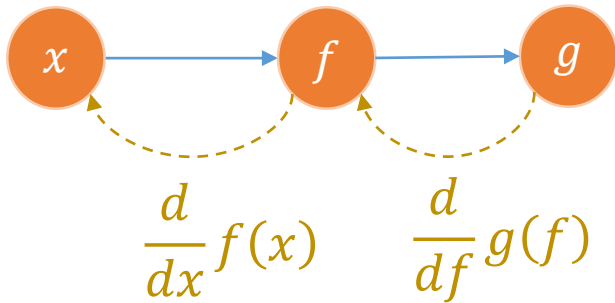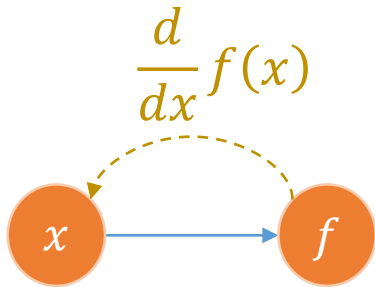$$\frac{d}{dx} f(x, y)$$

19

# Outline

➢ **Review on Optimization**

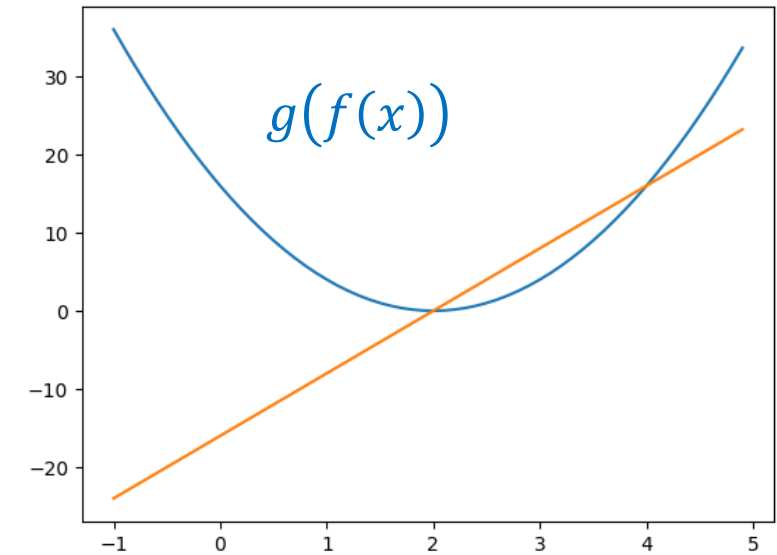➢ **Partial Gradient and Chain Rules**
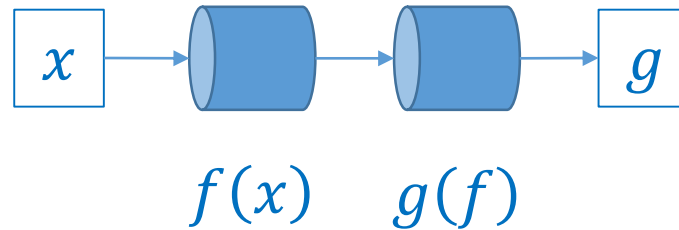
➢ **Finding a Line**

➢ **Linear Regression**

# Gradient-based Optimization
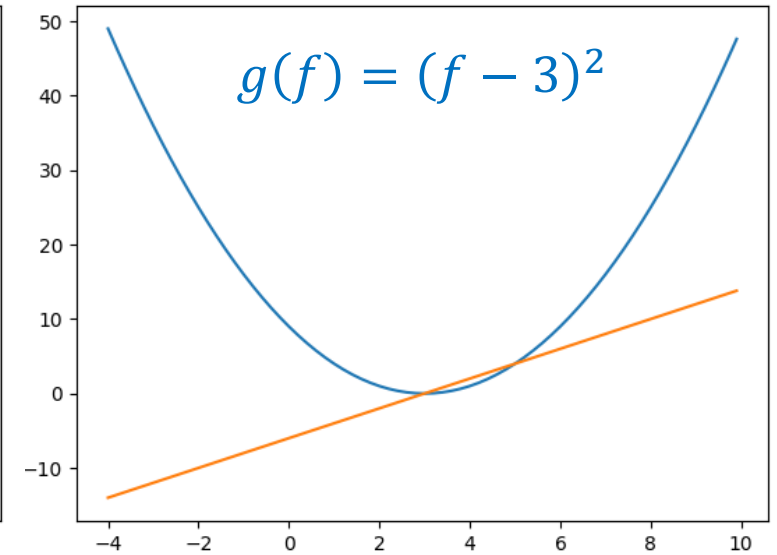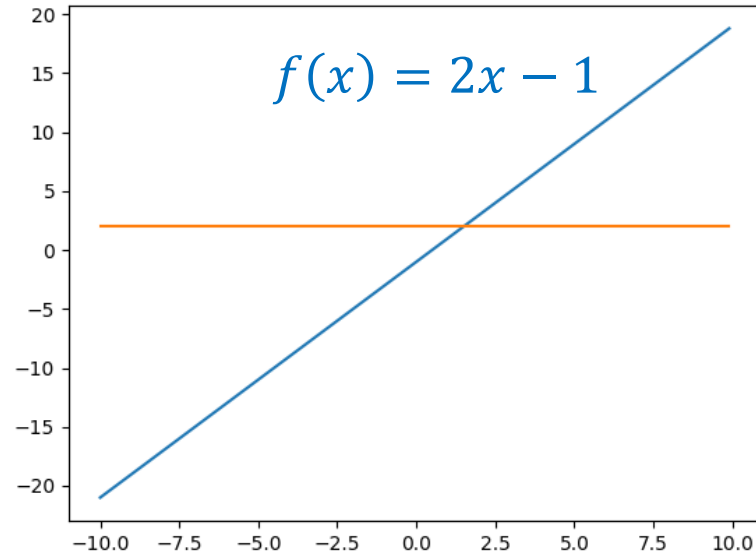
❖ **For composite function**

$$\frac{d}{dx}f(x)$$

$$x \rightarrow f$$

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

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

$$x \rightarrow f \rightarrow g$$

$$\frac{d}{dx}f(x) \qquad \frac{d}{df}g(f)$$

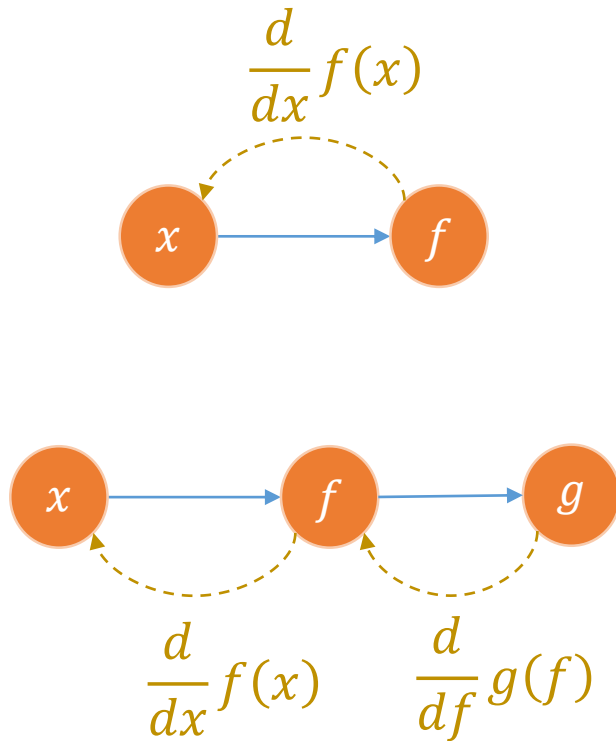$$x \rightarrow f(x) \rightarrow g(f) \rightarrow g$$

$$g(f(x))$$

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

# Gradient-based Optimization

❖ **For composite function**



$$\frac{d}{dx}f(x)$$

$$\frac{d}{dx}f(x) \qquad \frac{d}{df}g(f)$$

$$\frac{d}{dx}g\big(f(x)\big) = \left[\frac{d}{df}g(f)\right] * \left[\frac{d}{dx}f(x)\right]$$
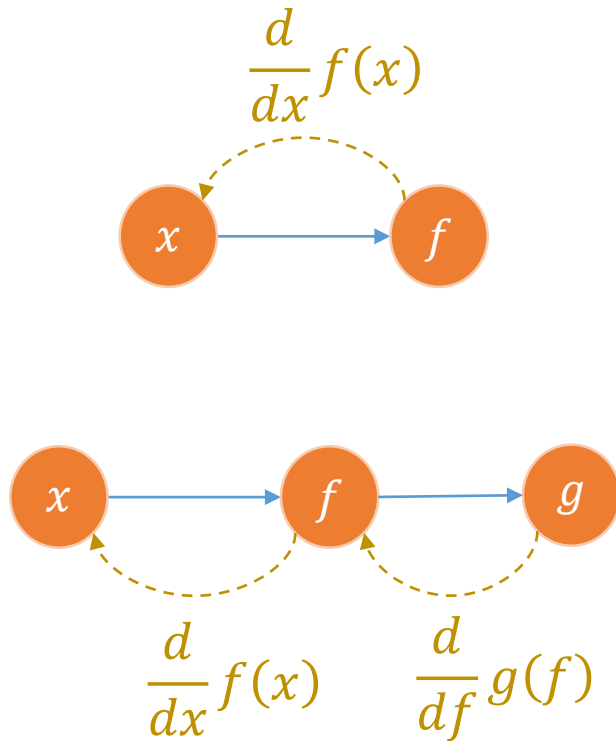
$$f(x) = 2x - 1$$
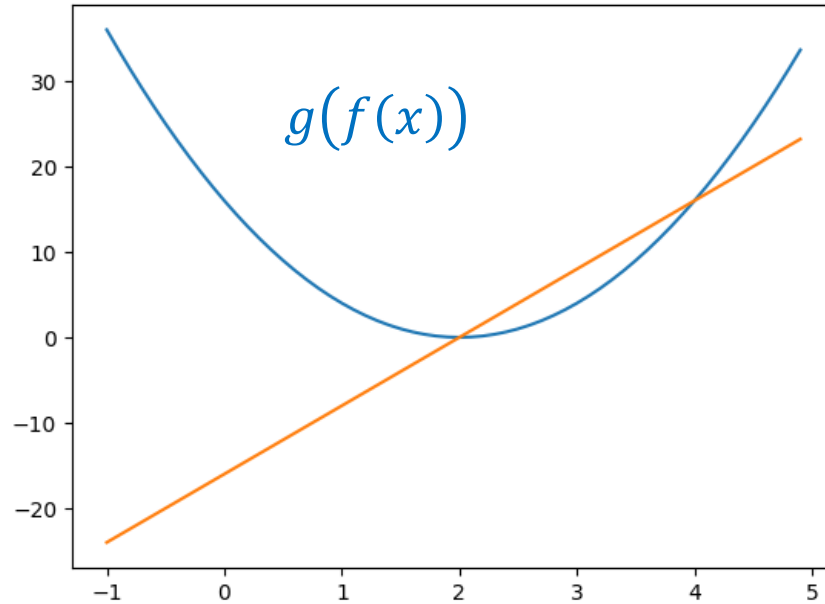
$$g(f) = (f - 3)^2$$

➡

$$g(x) = (2x - 1 - 3)^2$$

$$= (2x - 4)^2$$

➡

$$g'(x) = 4(2x - 4)$$

$$= 8x - 16$$

22

# Gradient-based Optimization

❖ **For composite function and chain rule**

$$\frac{d}{dx}f(x)$$



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

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

➡

$$f'(x) = 2$$

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

$$\frac{d}{dx}f(x) \qquad \frac{d}{df}g(f)$$

➡

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

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

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

$$= 8x - 16$$

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

**Implementation**

$$g(f(x))$$

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

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

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

```python
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
```

```python
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}')
```
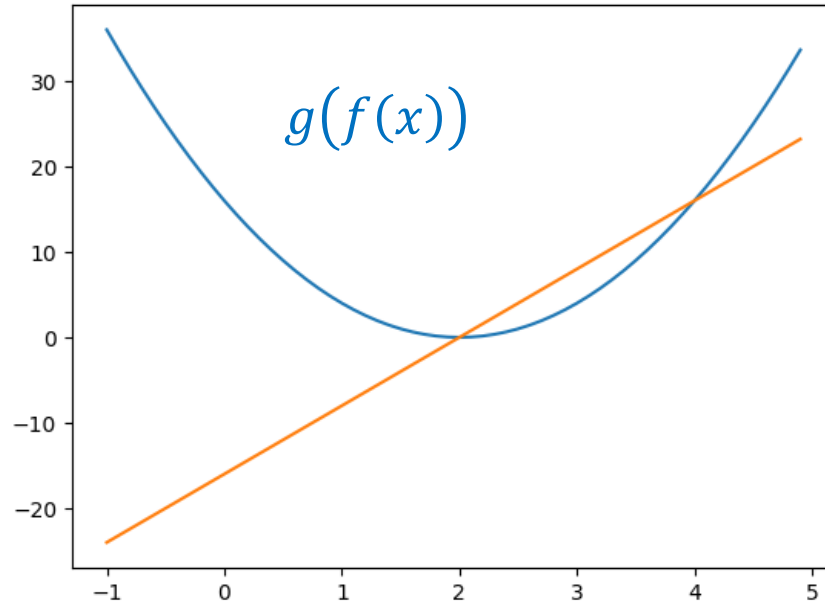
```
old_loss: 13456
new_loss: 538.2399999999994
```

**Implementation**


$$g(f(x))$$

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

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

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

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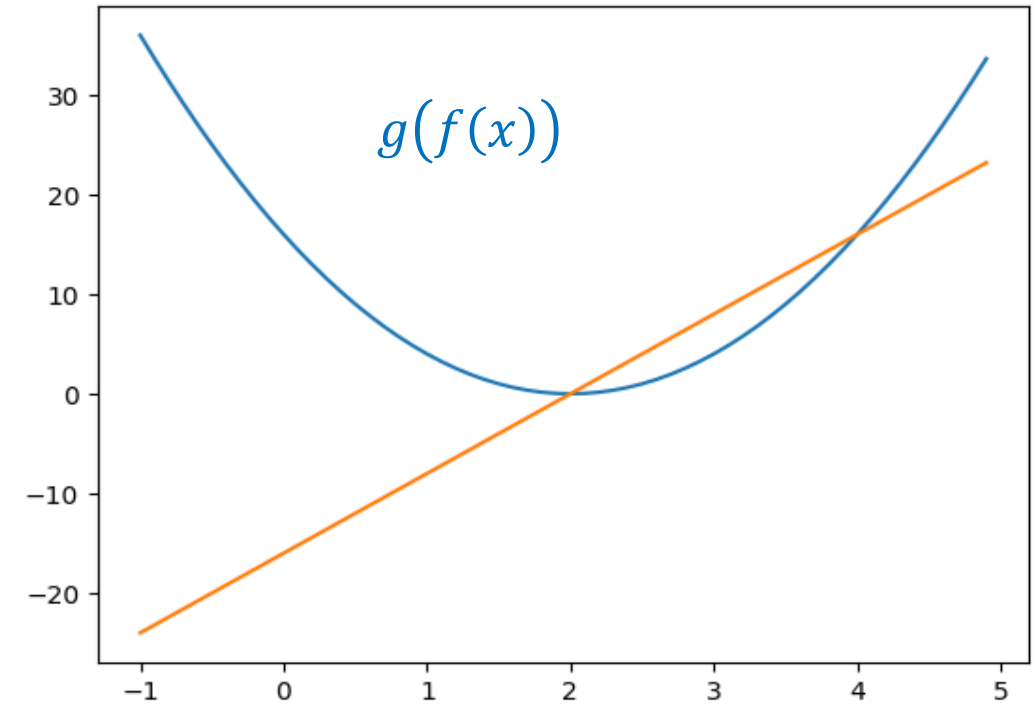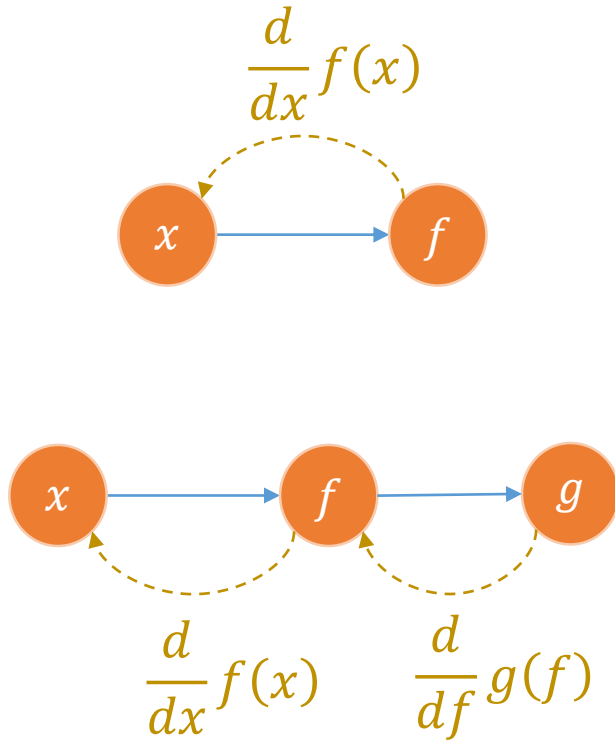
```python
1   import random
2
3   # parameters
4   num_steps = 5
5   lr = 0.1
6
7   # set x randomly
8   x = random.randint(-100, 100)
9
10  for _ in range(num_steps):
11      # Logging
12      loss = gf(fx(x))
13
14      # compute derivative
15      dg_dx_value = dg_dx(x)
16
17      # update
18    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

# Gradient-based Optimization

❖ **For composite function**

$$\frac{d}{dx}f(x)$$

$$x \rightarrow f$$

$$x \rightarrow f \rightarrow g$$

$$\frac{d}{dx}f(x) \qquad \frac{d}{df}g(f)$$

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

$$g(f(x))$$

- Given a function g(f(x)), find optimal $x_{opt}$ so that $f(x_{opt})$ is minimum

- After an update, $g(f(x_{new})) \leq g(f(x_{old}))$

# Gradient-based Optimization
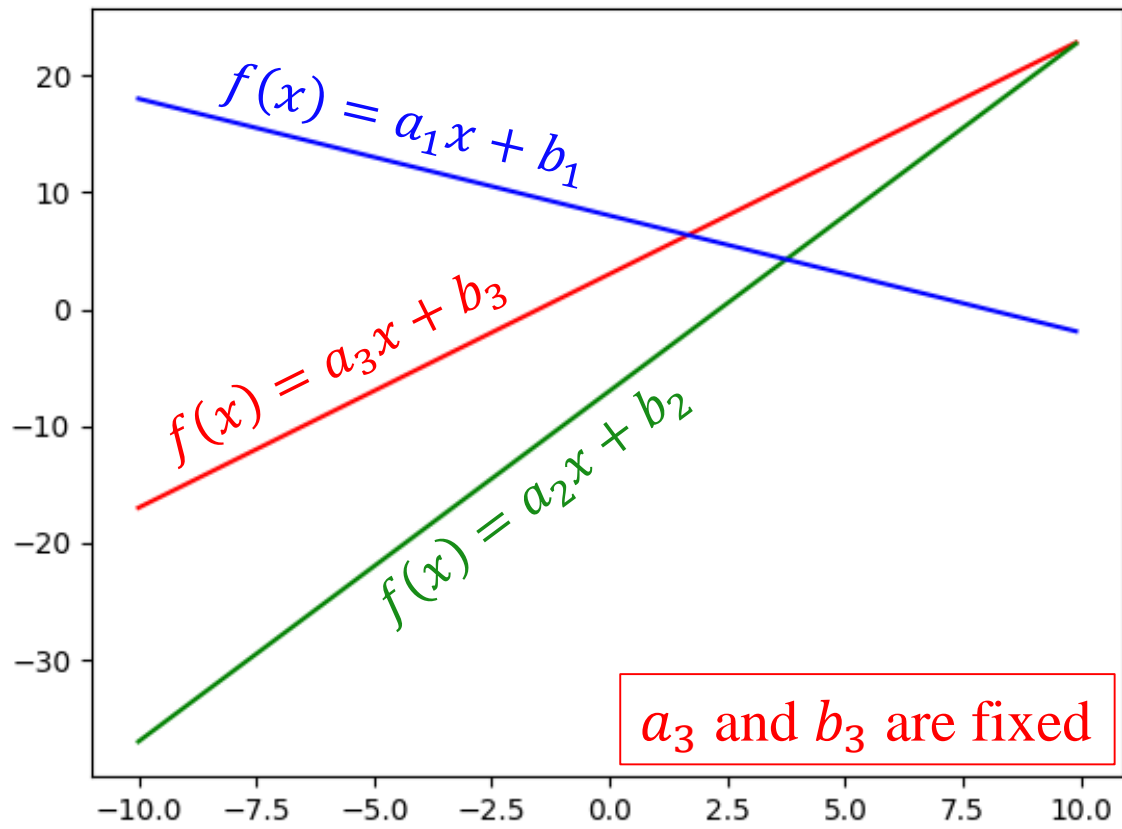
❖ **Another context**

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



$f(x) = a_1 x + b_1$

$f(x) = a_3 x + b_3$

$f(x) = a_2 x + b_2$

$a_3$ and $b_3$ are fixed

❖ **Constraints**

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

x → The red line (Black box) → y

How to measure the distance between the red line and an predicted line?

Initialize a, b

Compute partial gradient at a, b
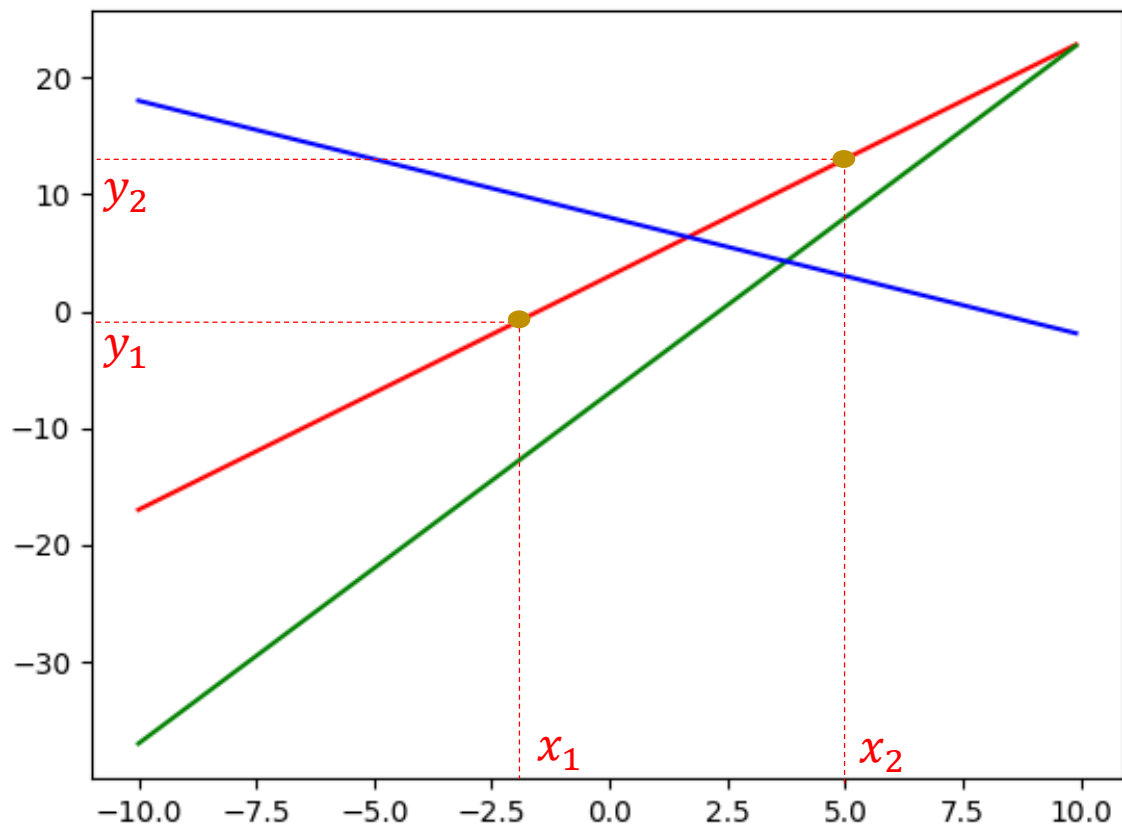
Move a, b opposite to db, db

$\eta = 0.1$

# Gradient-based Optimization

## ❖ Another context

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



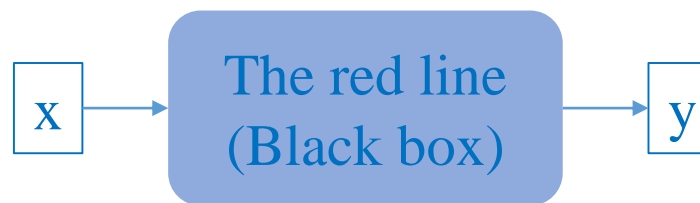## ❖ Constraints

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

x → The red line (Black box) → y

$(x_1 = -2, y_1 = -1)$

$(x_2 = 5, y_2 = 13)$

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

Initialize a, b

Compute partial gradient at a, b

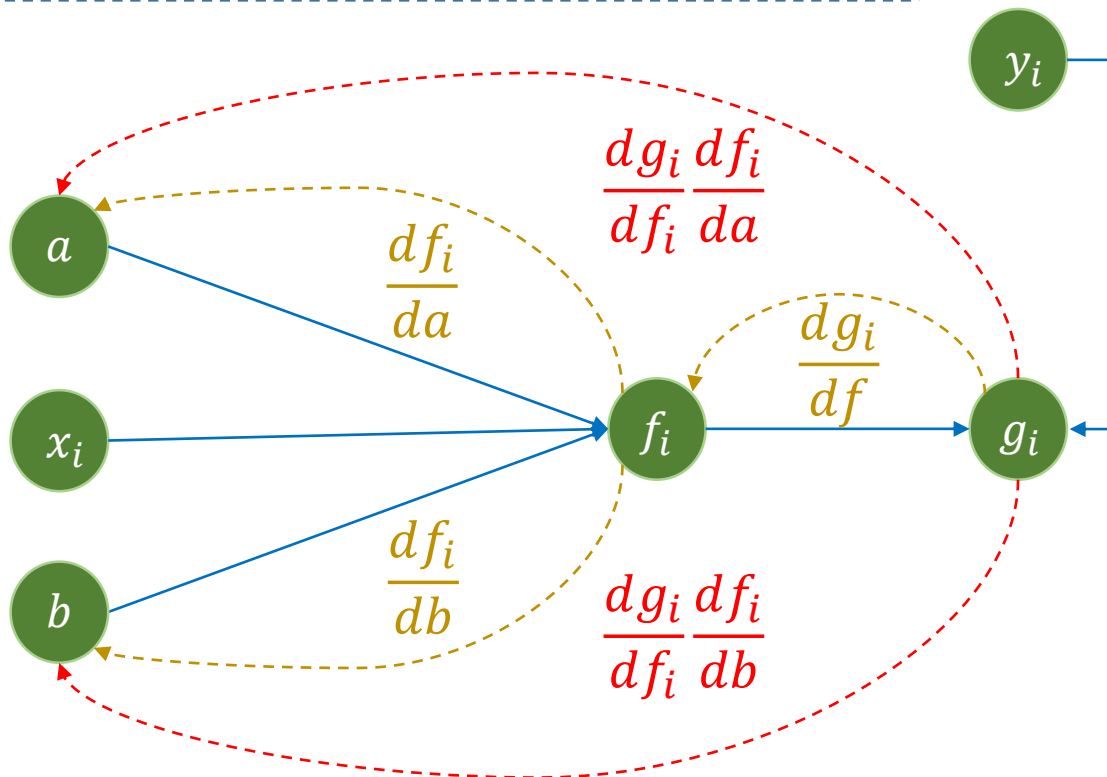Move a, b opposite to db, db

$\eta = 0.1$

# Gradient-based Optimization

❖ **Equations for partial gradients**

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

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



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

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

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

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

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

29

# ❖ Optimization for a composite function

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

$$f(x_i) = ax_i + b \qquad (x_1=1, y_1=5)$$
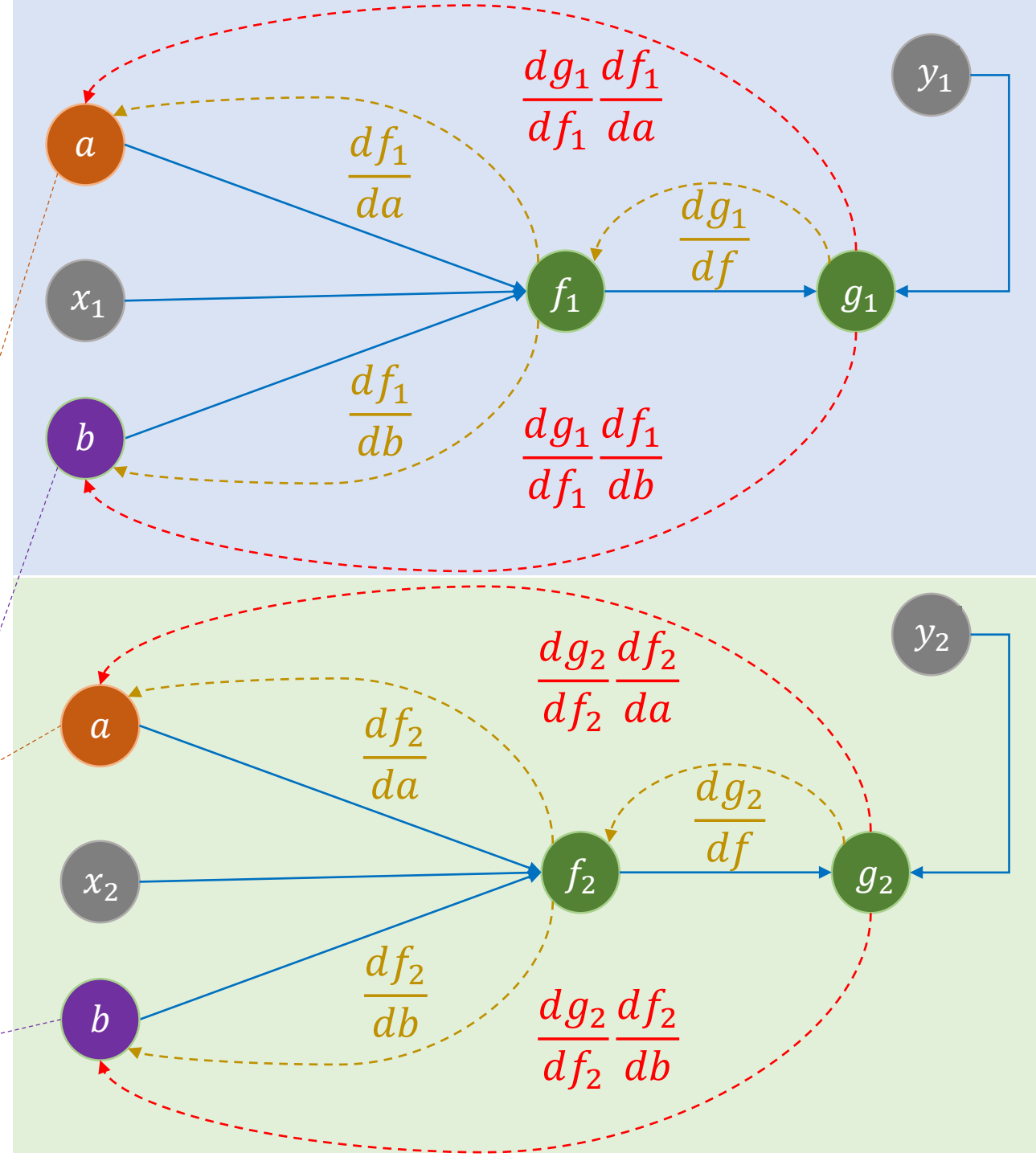
$$g(f) = (f - y_i)^2 \qquad (x_2=2, y_2=7)$$

## Partial derivative functions

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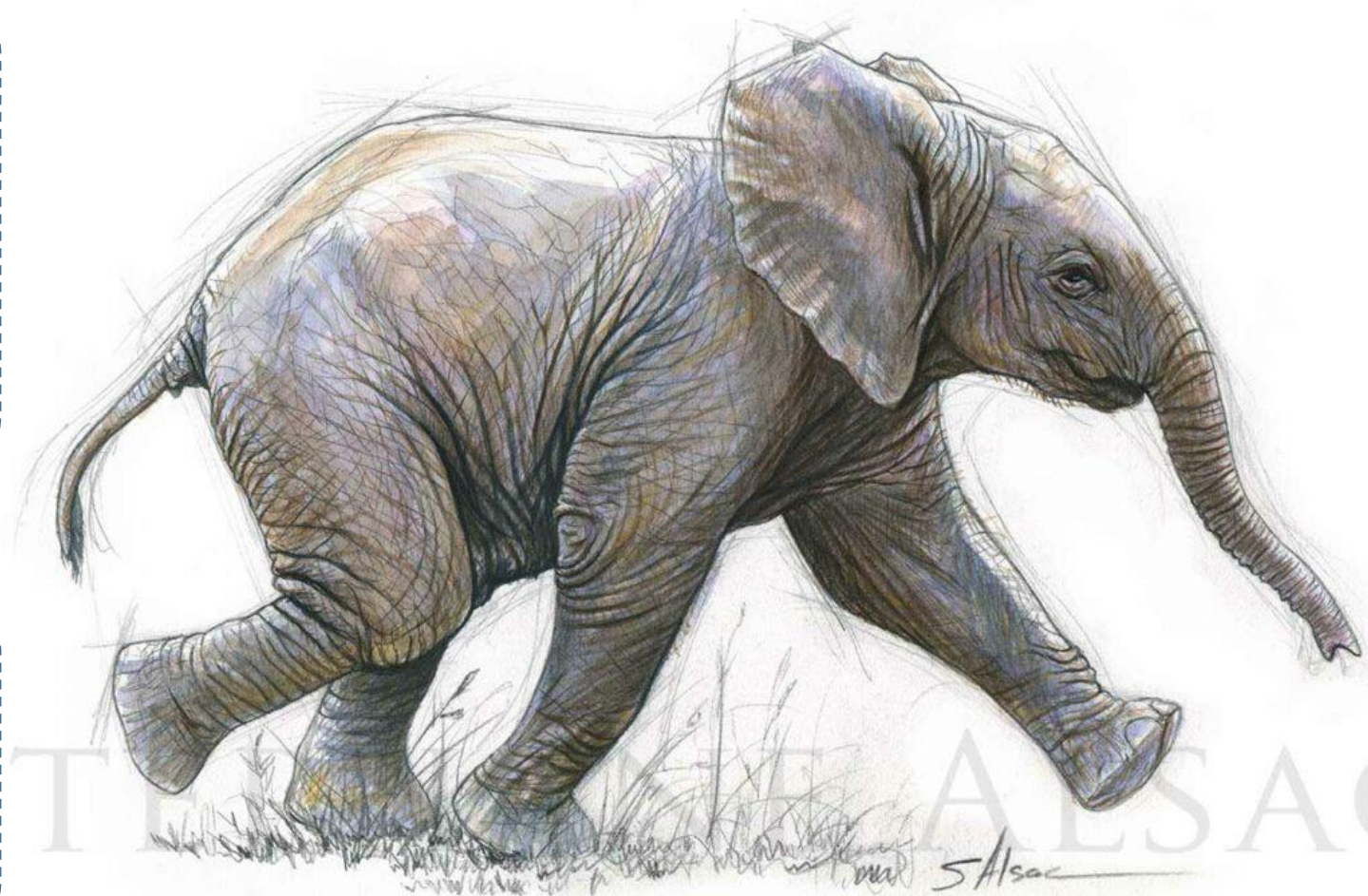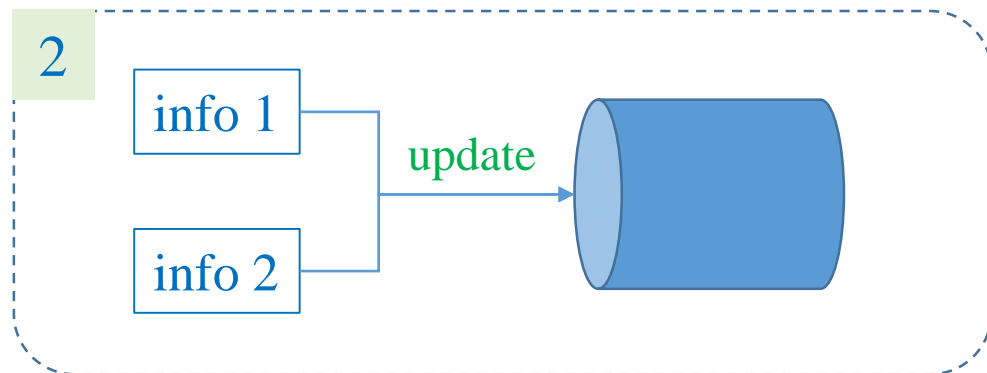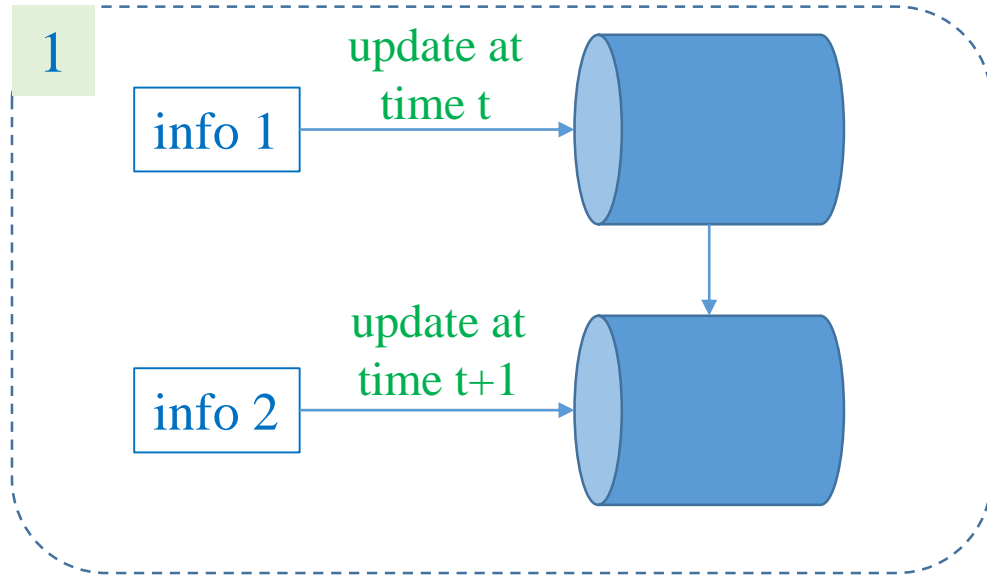
$$\frac{dg}{db} = \frac{dg}{df}\frac{df}{db} = 2(f-y)$$

$$\sum_i \frac{dg_i}{da} = \frac{dg_1}{df_1}\frac{df_1}{da} + \frac{dg_2}{df_2}\frac{df_2}{da}$$

$$\sum_i \frac{dg_i}{db} = \frac{dg_1}{df_1}\frac{df_1}{db} + \frac{dg_2}{df_2}\frac{df_2}{db}$$
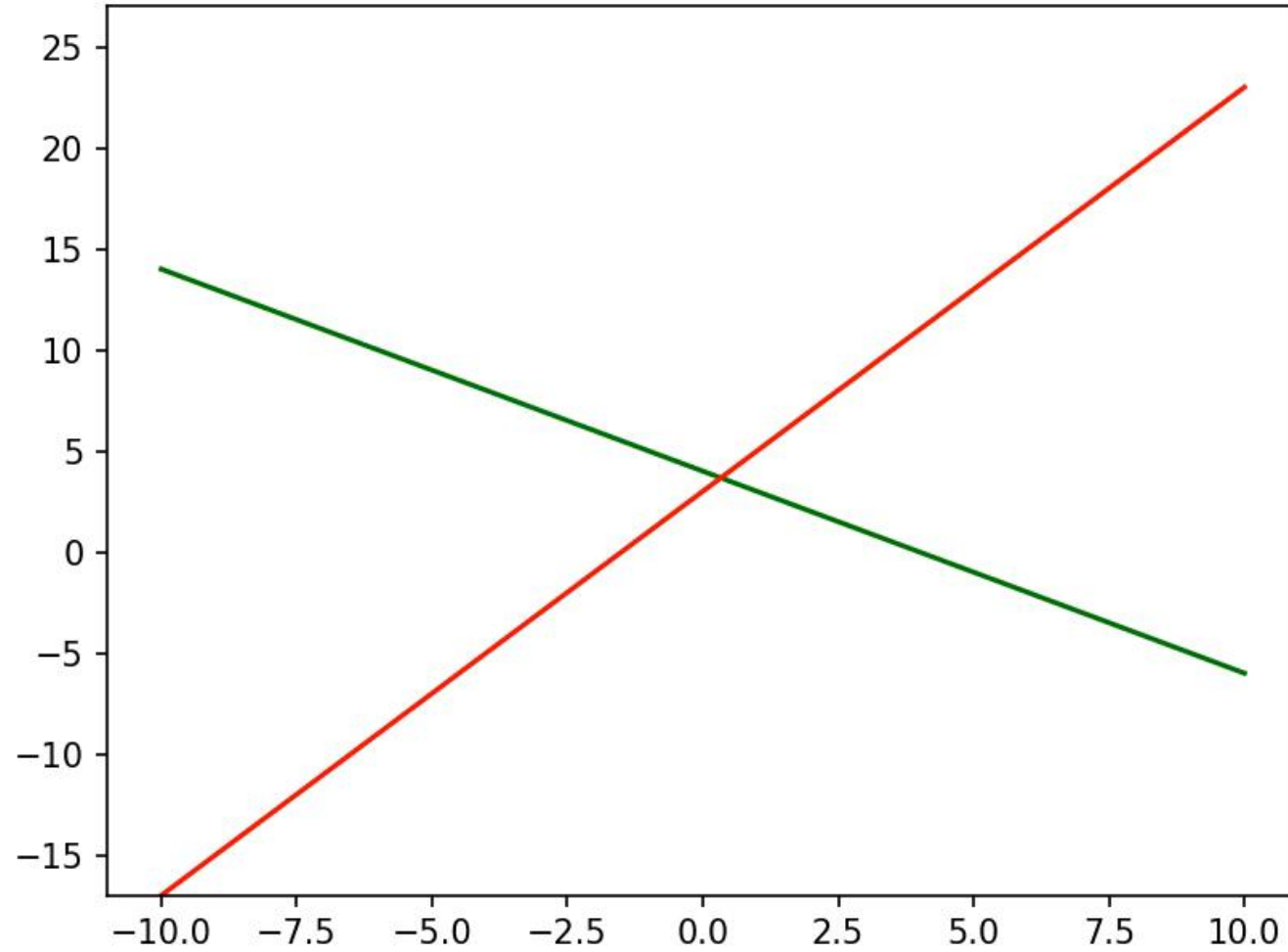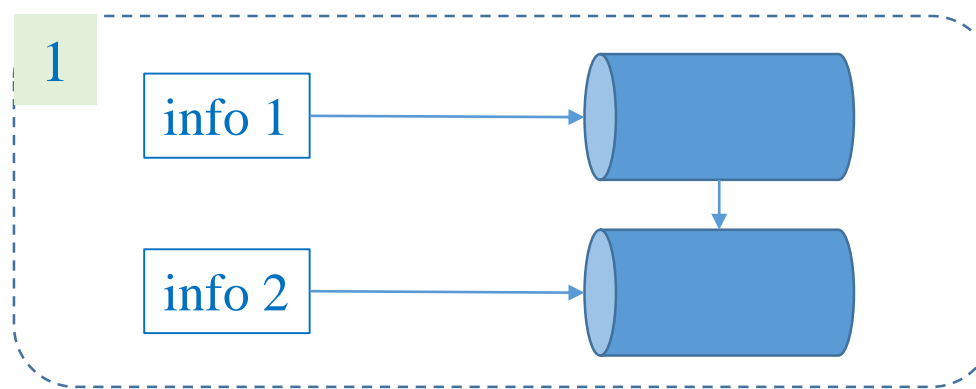
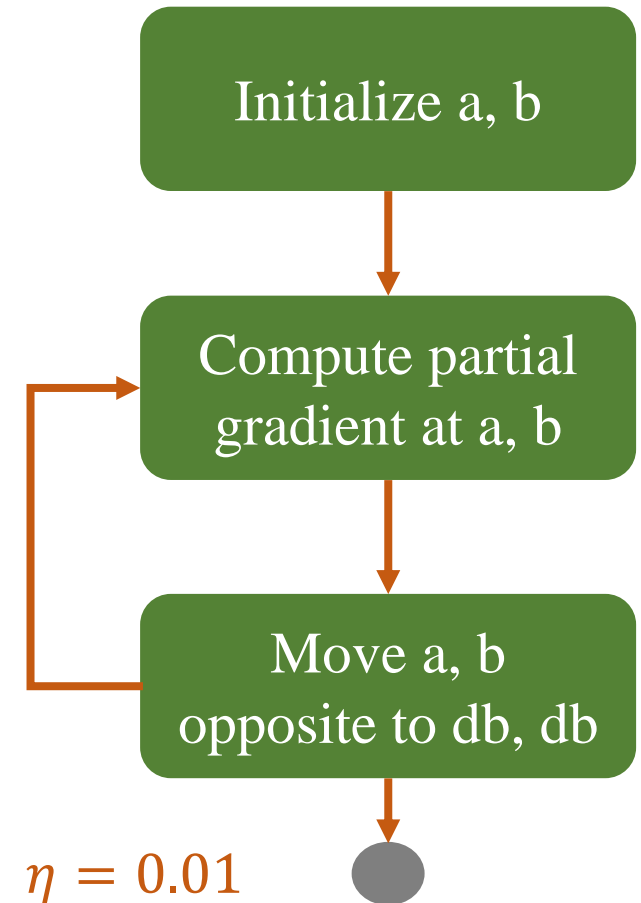# Gradient-based Optimization

❖ **How to use gradient information**

1

info 1 → update at time t → [cylinder]

↓

info 2 → update at time t+1 → [cylinder]

2

info 1 ⌐
          ├ update → [cylinder]
info 2 ⌐

```python
import random

# predict function
def predict_func(x, a, b):
    return a*x + b

# parameters
num_steps = 100
lr = 0.01

# given data
data = [[-2, -1],
        [5, 13]]

# 1. set a, b randomly
a = random.random()*10.0 - 5.0
b = random.random()*10.0 - 5.0
```

```python
for i in range(num_steps):
    for sample in data:
        x_value, y_value = sample

        # compute predicted_y
        predicted_y = predict_func(x_value, a, b)


        # compute g
        g_value = (predicted_y - y_value)**2


        # compute partial gradients for a and b
        dg_da = 2*x_value*(predicted_y - y_value)
        dg_db = 2*(predicted_y - y_value)

        # update
        a = a - lr*dg_da
        b = b - lr*dg_db
```

# Summary

1

info 1 → [cylinder]

info 2 → [cylinder]

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

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



Initialize a, b

↓

Compute partial gradient at a, b

↓

Move a, b opposite to db, db

$\eta = 0.01$

33

```python
import random

# predict function
def predicted_func(x, a, b):
    return a*x + b


# parameters
num_steps = 100
lr = 0.01


# given data
x1 = -2
y1 = -1

x2 = 5
y2 = 13


# 1. set a, b randomly
a = random.random()*10.0 - 5.0
b = random.random()*10.0 - 5.0
```
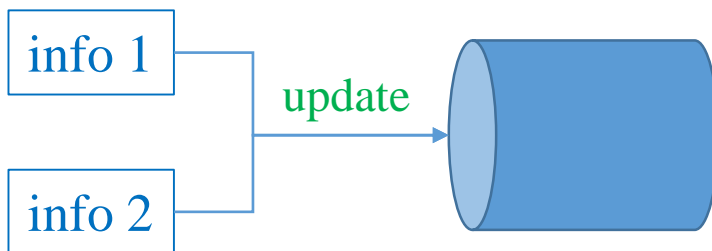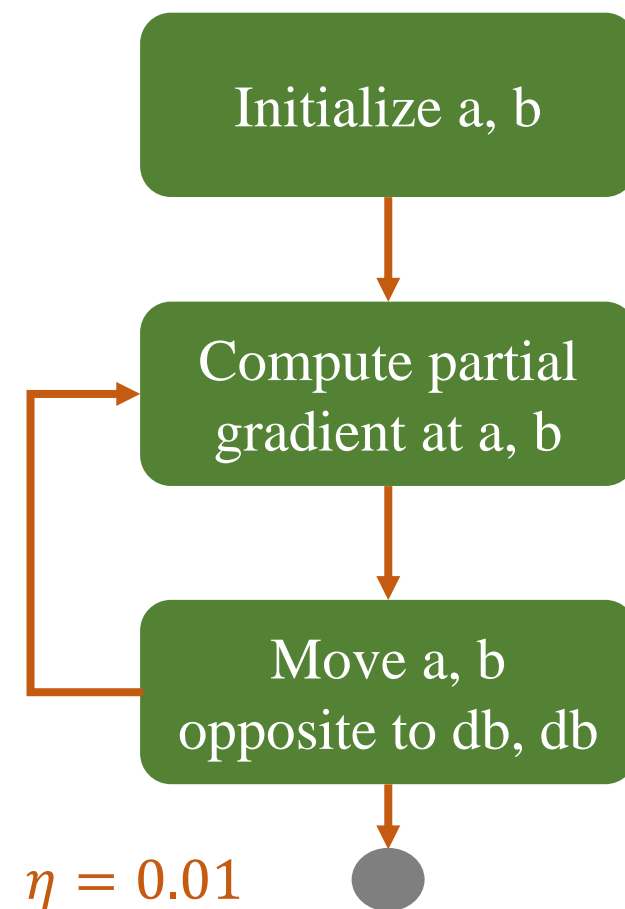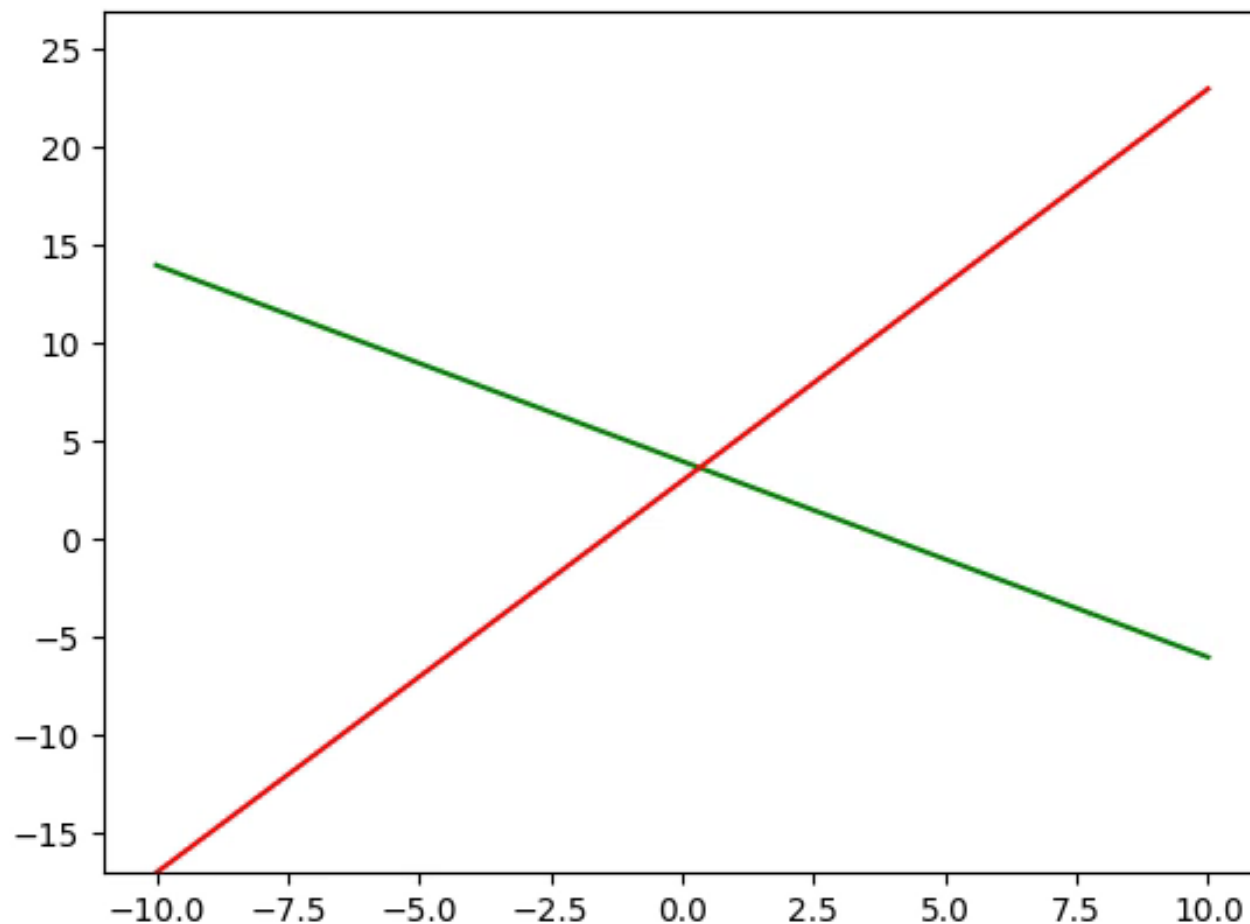
```python
for i in range(num_steps):
    # 2. compute predicted_y1 and predicted_y2
    predicted_y1 = predicted_func(x1, a, b)
    predicted_y2 = predicted_func(x2, a, b)

    # 3. compute g
    g_value_1 = (predicted_y1 - y1)**2
    g_value_2 = (predicted_y2 - y2)**2

    # Logging
    # ...

    # 4. compute partial gradients for a and b
    dg_da = 2*x1*(predicted_y1 - y1) + 2*x2*(predicted_y2 - y2)
    dg_db = 2*(predicted_y1 - y1) + 2*(predicted_y2 - y2)

    # 5. update
    a = a - lr*dg_da
    b = b - lr*dg_db
```
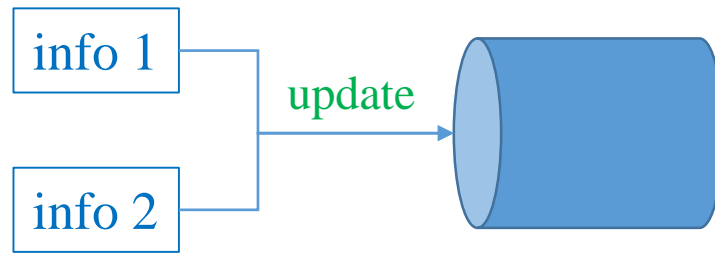
# Summary

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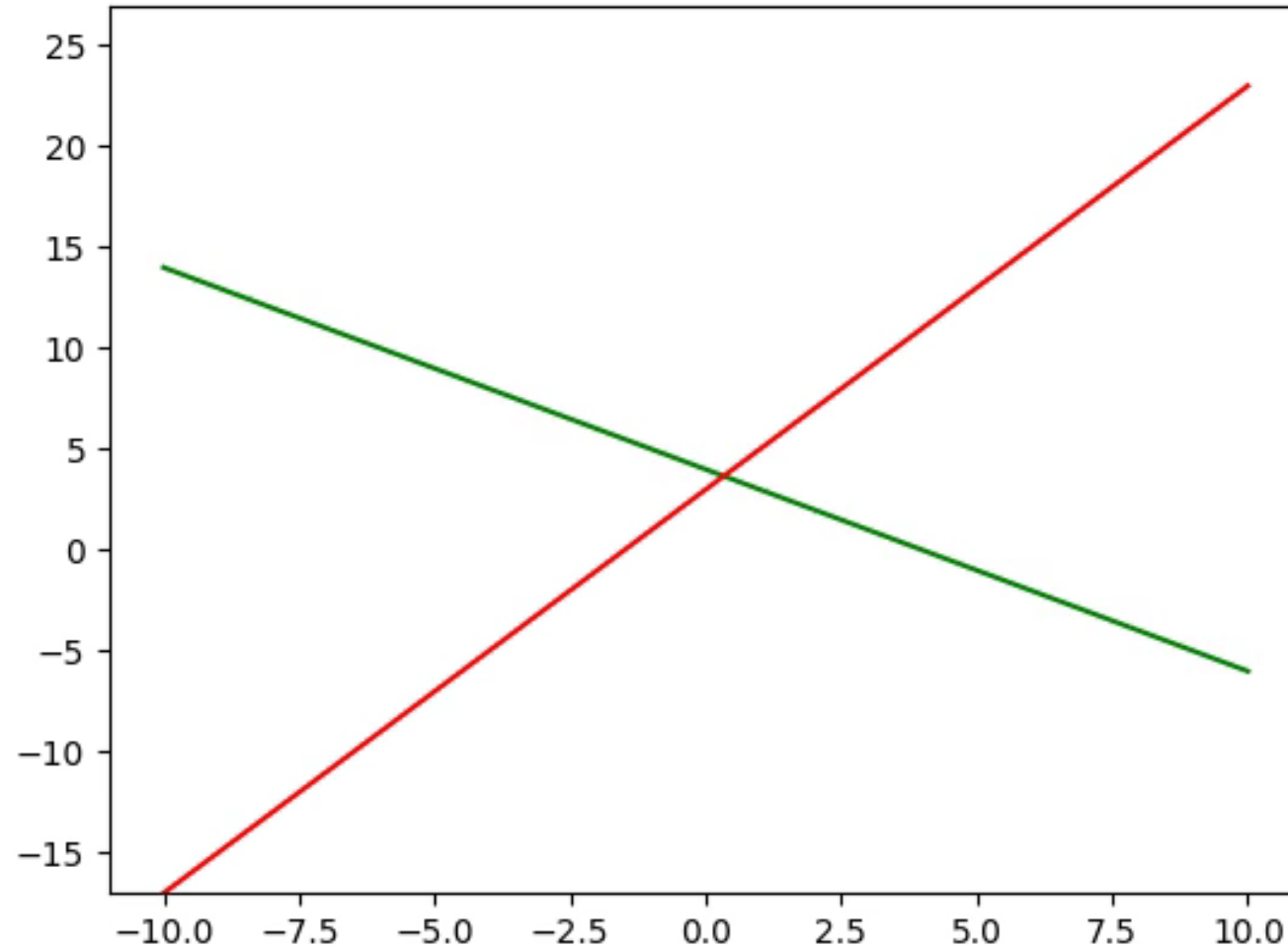
$$\frac{dg}{db} = \frac{dg}{df}\frac{df}{db} = 2(f-y)$$



Initialize a, b

Compute partial gradient at a, b

Move a, b opposite to db, db

$\eta = 0.01$

35

# 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 to db, db

$\eta = 0.001$

# List-based Implementation

❖ **Vectorization**

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

$$= ax + b1$$

$$= \begin{bmatrix} x \\ 1 \end{bmatrix} \cdot \begin{bmatrix} a \\ b \end{bmatrix}$$

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

```
1   # implement 1: naive approach
2   def predict_w1(x, a, b):
3       return a*x + b
4
5   # test
6   print(predict_w1(2, 1, 3))

5
```
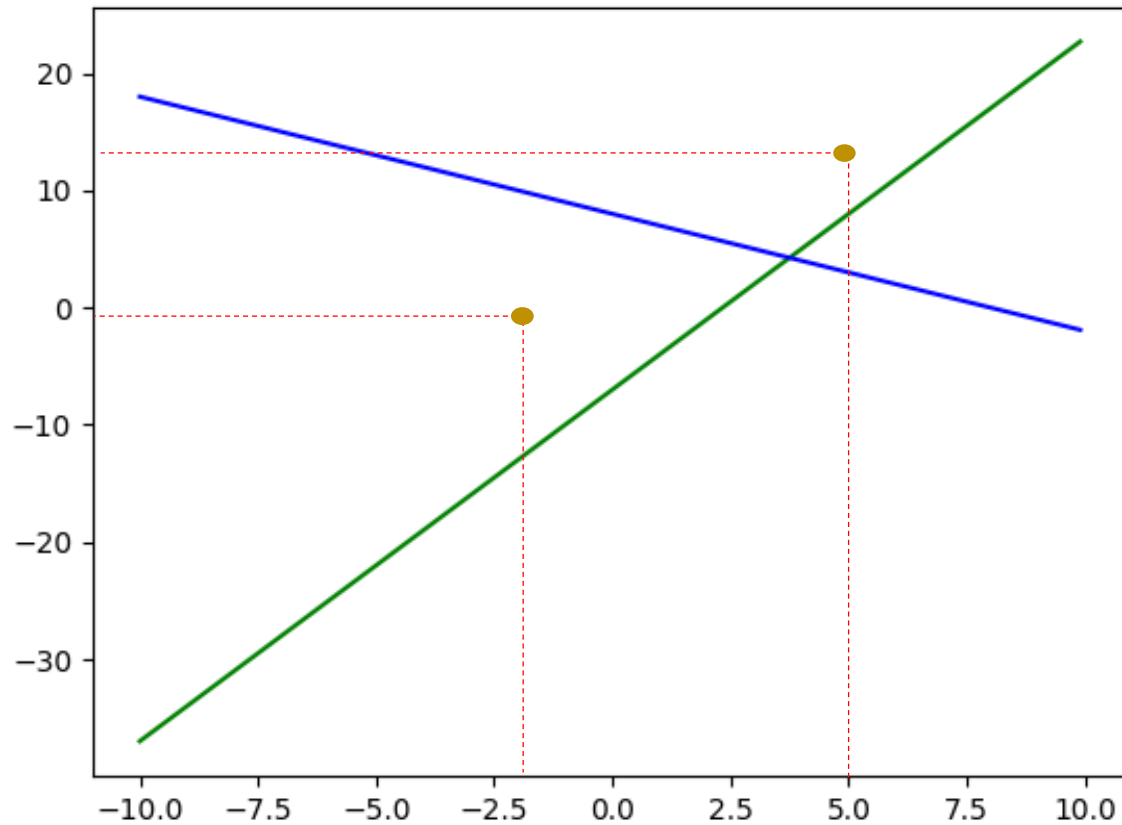
```
1   # implement 2: using list
2   def predict_w2(data, weights):
3       return sum([d*w for d, w in zip(data, weights)])
4
5   # test
6   print(predict_w2([2, 1], [1, 3]))

5
```
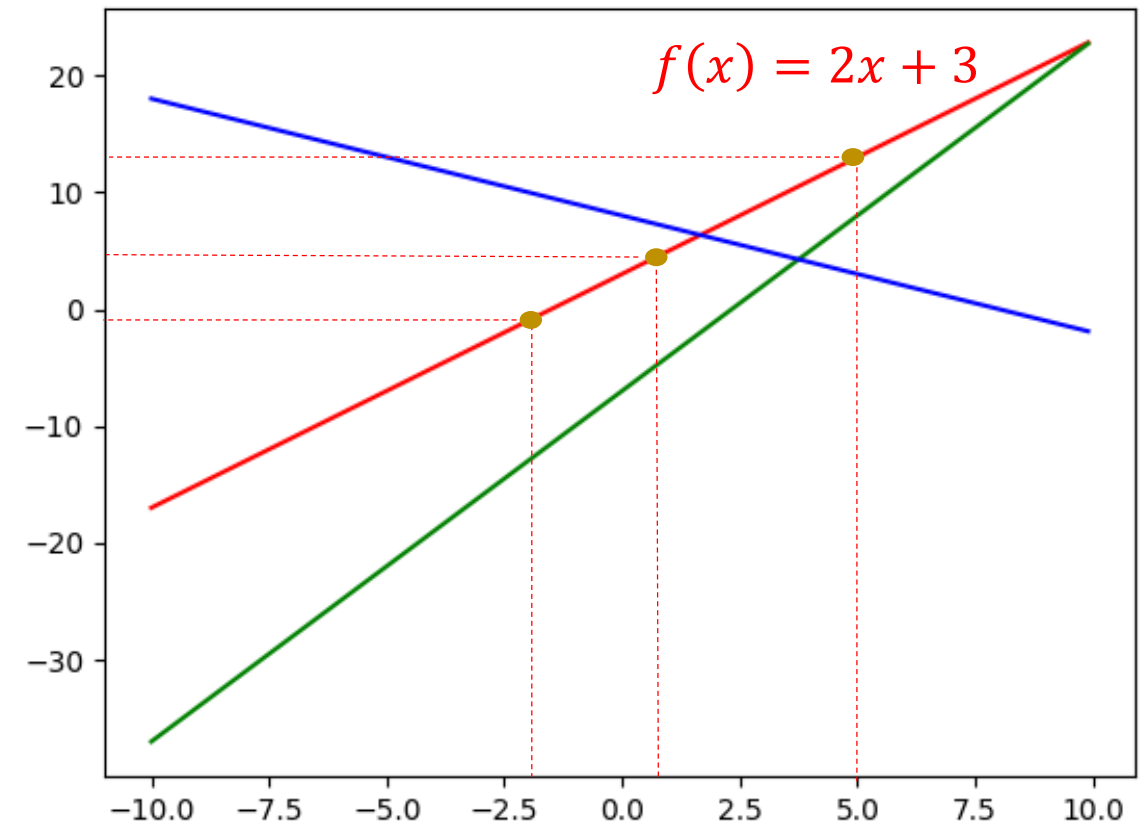
37

# Discussion

$(x_1 = -2, y_1 = -1)$

$(x_2 = 5, y_2 = 13)$

**Have one**

**more sample**

$(x_2 = 1, y_2 = 3)$

**Remove the red line**



$f(x) = 2x + 3$

# Discussion

❖ **What about the given samples?**



Line 1: go through two points

Line 2: smallest summation of distances

Line 3: go through one point