

Momentum

Gradient Descent with Momentum

Nesterov Momentum

Nguyen Quoc Thai

CONTENT

(1) – Multilayer Perceptron

(2) – Gradient Descent with Momentum

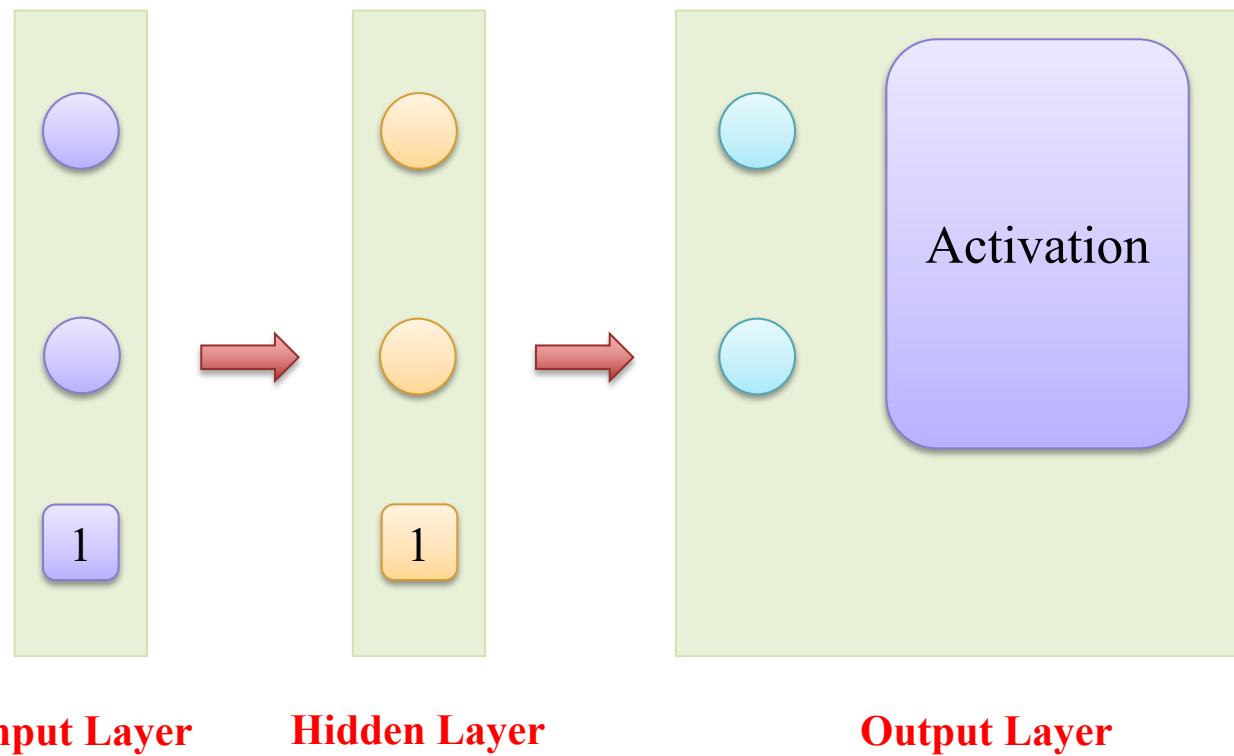
(3) – Nesterov Momentum

1 – Multilayer Perceptron

!

Multilayer Perceptron

#parameters: 12

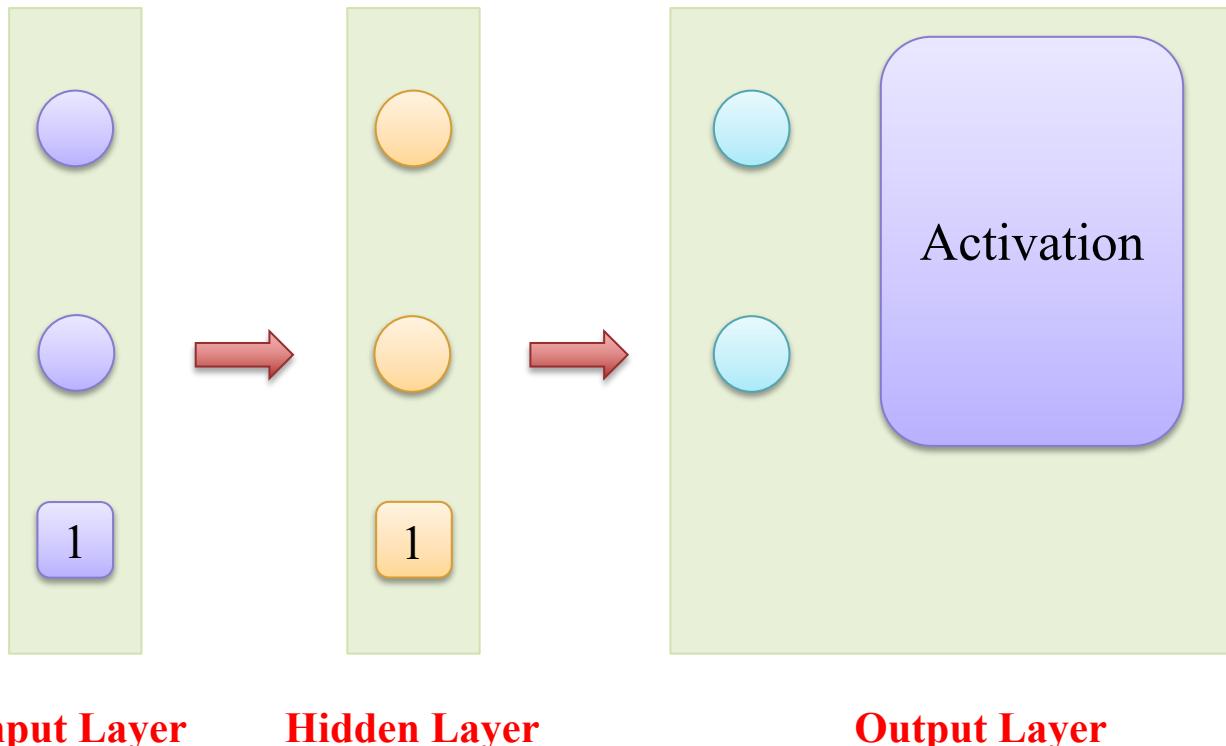


1 – Multilayer Perceptron



Multilayer Perceptron

#parameters: 12



```
model = nn.Sequential(  
    nn.Linear(2, 2),  
    nn.Linear(2, 2),  
    nn.Sigmoid()  
)
```

```
summary(model, (1, 2))
```

Layer (type)	Output Shape	Param #
Linear-1	[-1, 1, 2]	6
Linear-2	[-1, 1, 2]	6
Sigmoid-3	[-1, 1, 2]	0

```
Total params: 12
```

```
Trainable params: 12
```

```
Non-trainable params: 0
```

```
Input size (MB): 0.00
```

```
Forward/backward pass size (MB): 0.00
```

```
Params size (MB): 0.00
```

```
Estimated Total Size (MB): 0.00
```

1 – Multilayer Perceptron

!

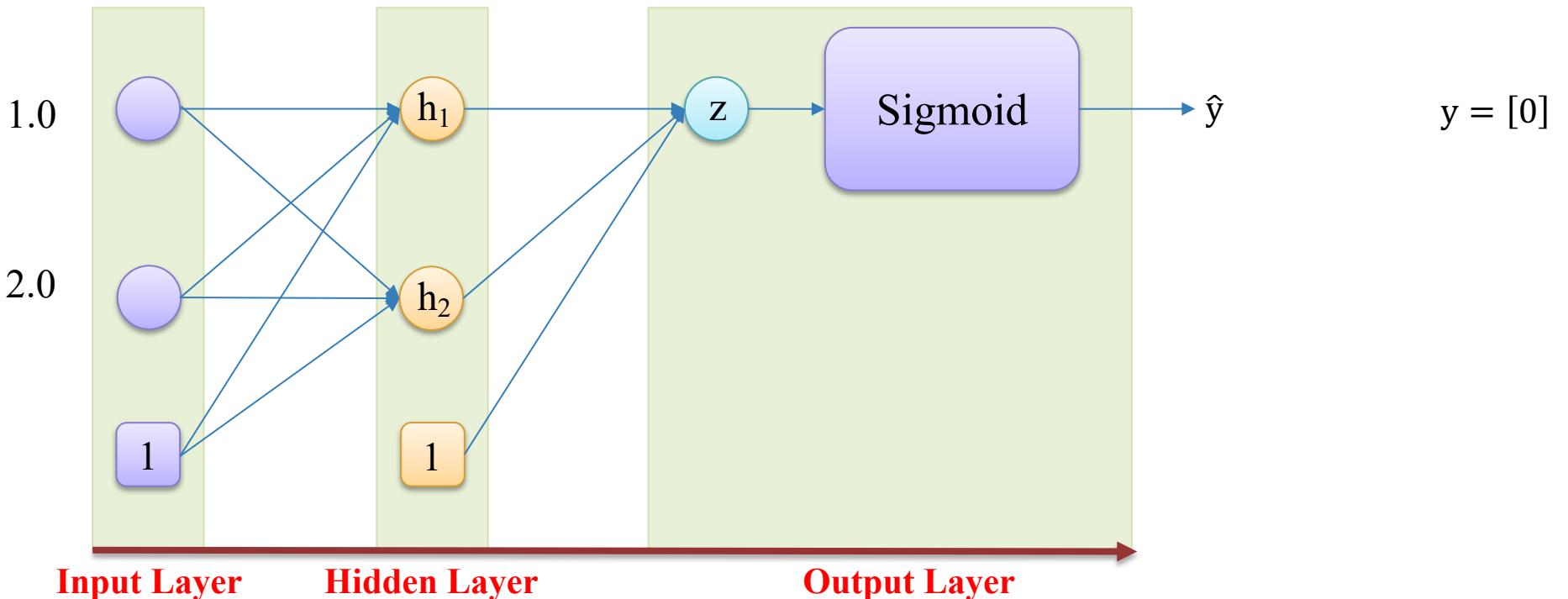
Forward

$$x = [1.0 \quad 2.0]$$

$$W_h = [W_{h1} \quad W_{h2}]$$

$$= \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix}$$

$$W_z = [W_z] = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$$



1 – Multilayer Perceptron

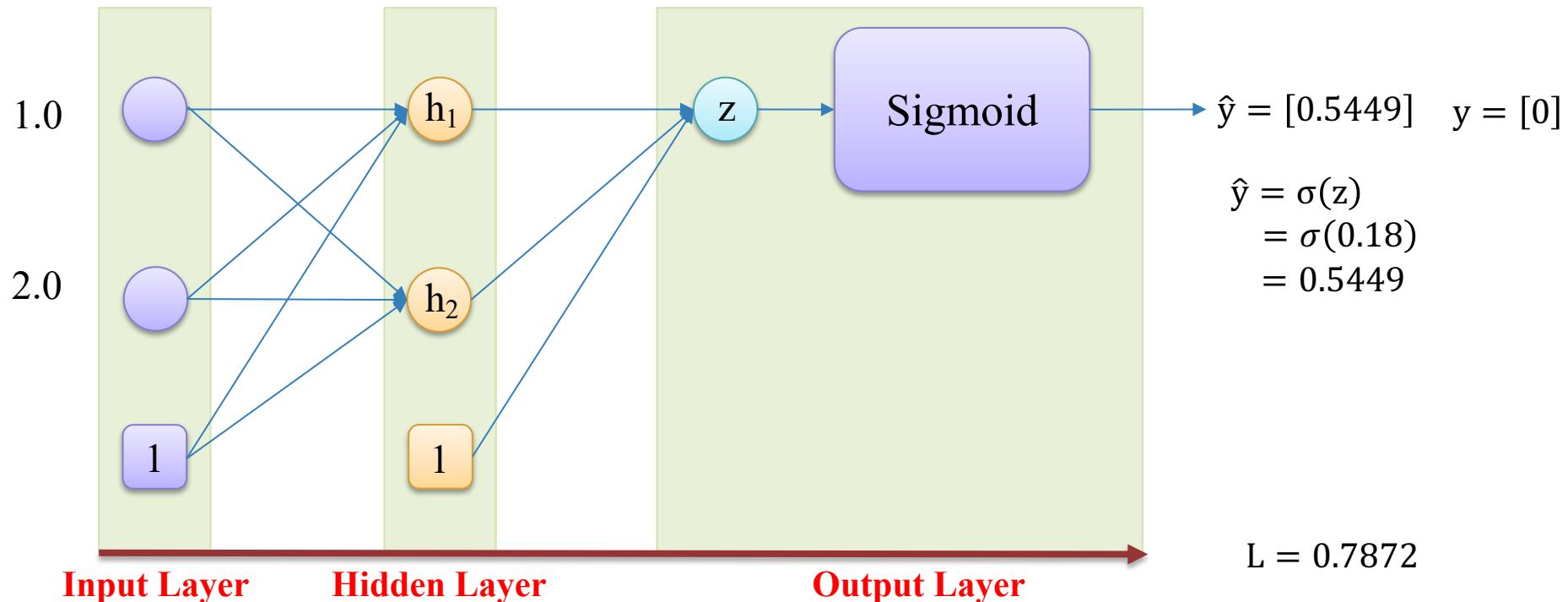
!

Forward

$$x = [1.0 \quad 2.0]$$

$$W_h = [W_{h1} \quad W_{h2}] \\ = \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix}$$

$$W_z = [W_z] = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$$



$$h = [1.0 \ x]W_h = [1.0 \quad 1.0 \quad 2.0] \begin{bmatrix} 0.1 & 0.1 \\ 0.1 & 0.1 \\ 0.1 & 0.1 \end{bmatrix} = [0.4 \quad 0.4]$$

$$z = [1.0 \ h]W_z = [1.0 \quad 0.4 \quad 0.4] \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix} = 0.18$$

1 – Multilayer Perceptron

!

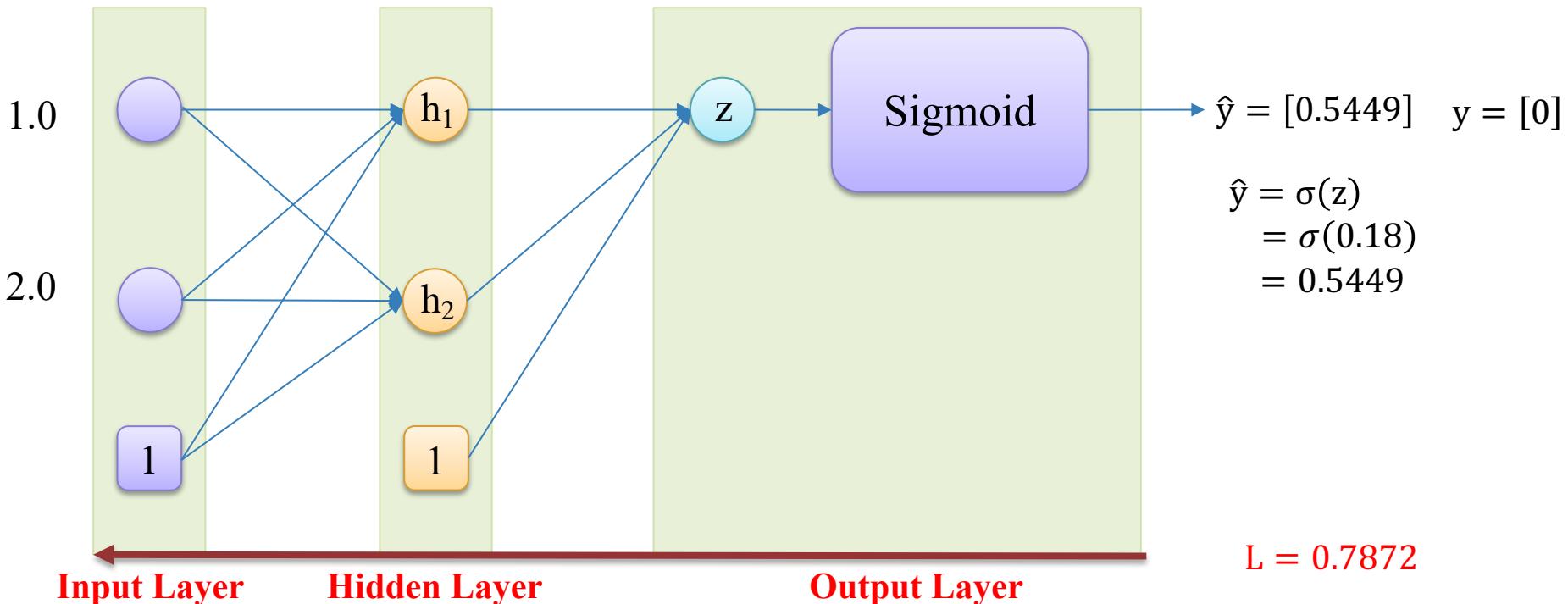
Backward

$$x = [1.0 \quad 2.0]$$

$$W_h = [W_{h1} \quad W_{h2}]$$

$$= \begin{bmatrix} 0.0946 & 0.0946 \\ 0.0946 & 0.0891 \\ 0.0946 & 0.0891 \end{bmatrix}$$

$$W_z = [W_z] = \begin{bmatrix} 0.0455 \\ 0.0782 \\ 0.0782 \end{bmatrix}$$



1 – Multilayer Perceptron



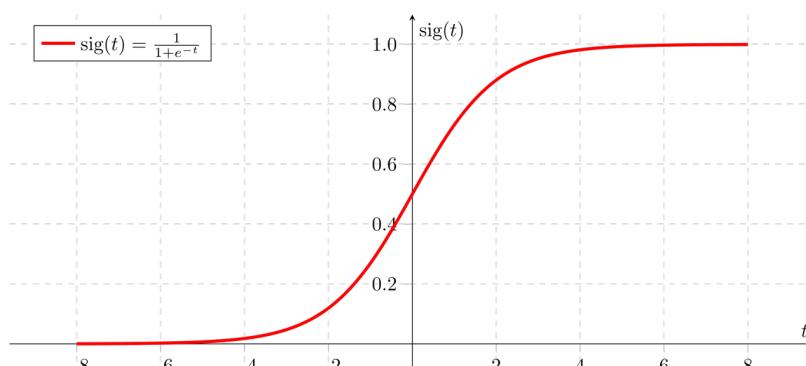
Activation

❖ Sigmoid Function

```
import torch.nn as nn

act = nn.Sigmoid()
input = torch.tensor([0.18, -0.18])
act(input)

tensor([0.5449, 0.4551])
```

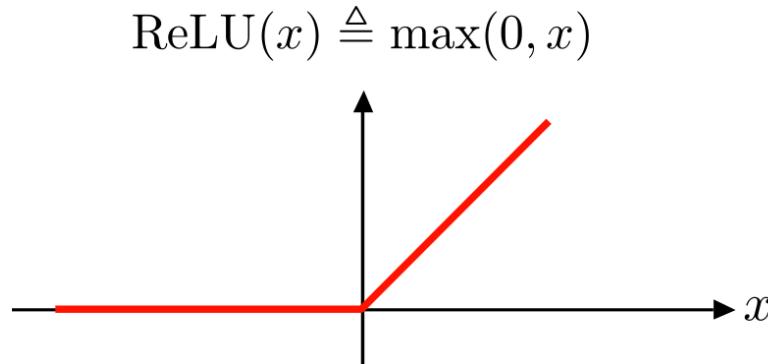


❖ ReLU Function

```
import torch.nn as nn

act = nn.ReLU()
input = torch.tensor([0.18, -0.18])
act(input)

tensor([0.1800, 0.0000])
```

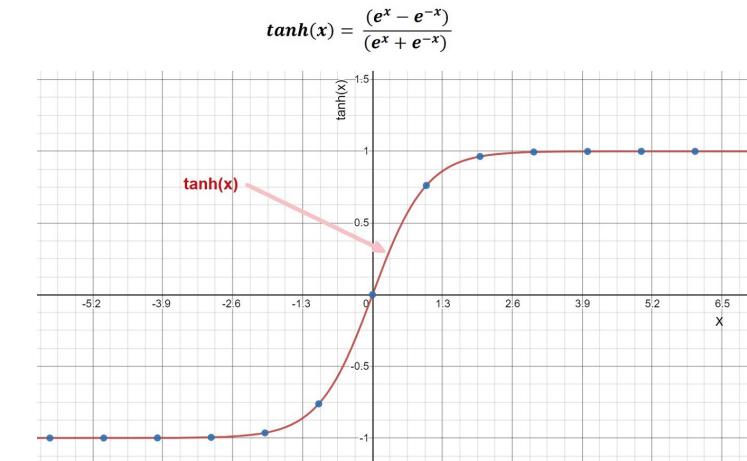


❖ Tanh Function

```
import torch.nn as nn

act = nn.Tanh()
input = torch.tensor([0.18, -0.18])
act(input)

tensor([ 0.1781, -0.1781])
```



1 – Multilayer Perceptron



Loss

❖ BCELoss()

```
import torch.nn as nn
loss_fn = nn.BCELoss()

y_pred
tensor([0.5449], grad_fn=<SigmoidBackward0>

y
tensor([0.])

loss = loss_fn(y_pred, y)
loss
tensor(0.7872, grad_fn=<BinaryCrossEntropyBackward0>)
```

❖ CrossEntropyLoss()

```
import torch.nn as nn
loss_fn = nn.CrossEntropyLoss()

y = torch.tensor(0)
y
tensor(0)

y_pred
tensor([0.0580, 0.4275], grad_fn=<AddBackward0>

loss_fn(y_pred, y)
tensor(0.8949, grad_fn=<NllLossBackward0>)
```

1 – Multilayer Perceptron



Optimizer

❖ SGD()

```
for layer in model.children():
    print(layer.state_dict())

OrderedDict([('weight', tensor([[0.1000, 0.1000],
                               [0.1000, 0.1000]])), ('bias', tensor([0.1000, 0.1000]))])
OrderedDict()
OrderedDict([('weight', tensor([[0.1000, 0.1000],
                               [0.1000, 0.1000]])), ('bias', tensor([0.1000, 0.1000]))])

learning_rate = 0.1
optimizer = optim.SGD(model.parameters(), learning_rate)

loss.backward()

optimizer.step()

for layer in model.children():
    print(layer.state_dict())

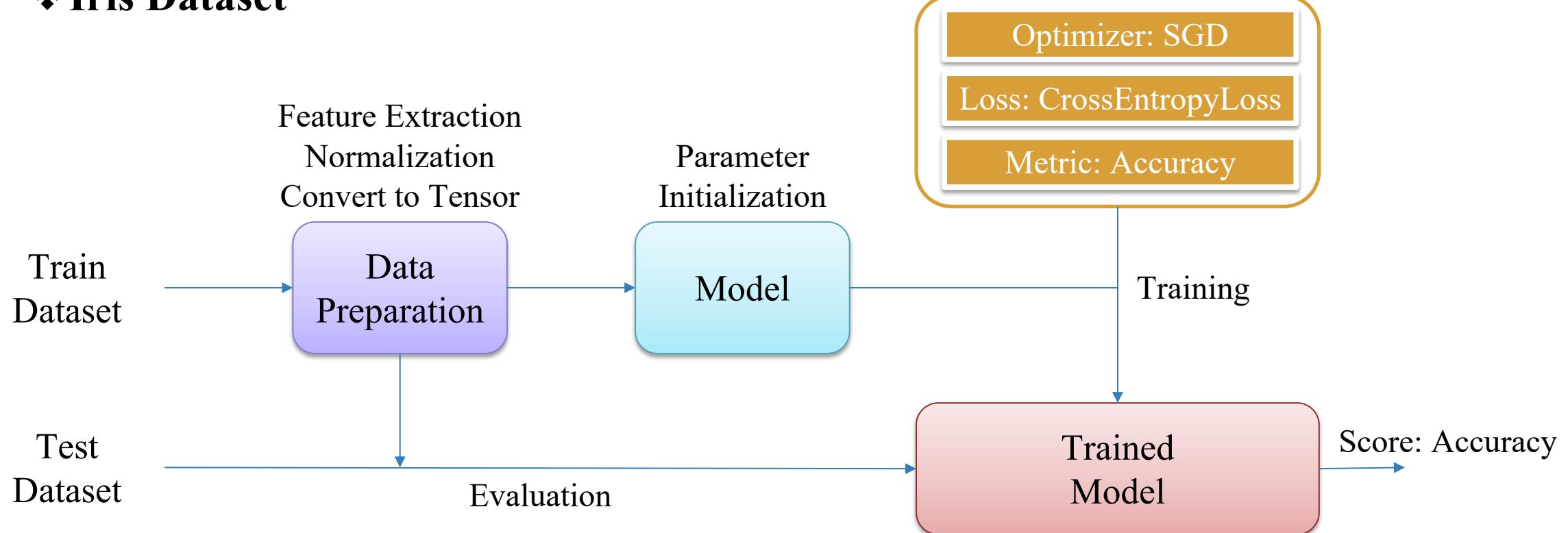
OrderedDict([('weight', tensor([[0.0946, 0.0891],
                               [0.0946, 0.0891]])), ('bias', tensor([0.0946, 0.0946]))])
OrderedDict([('weight', tensor([[0.0782, 0.0782]])), ('bias', tensor([0.0455]))])
OrderedDict()
```

1 – Multilayer Perceptron

!

Classification using Multilayer Perceptron

❖ Iris Dataset



2 – GD with Momentum

!

Gradient Descent

- 1) Pick a sample (x, y) from training data
- 2) Compute output \hat{y}

$$z = \boldsymbol{\theta}^T \mathbf{x}$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

- 3) Compute loss

$$L(\boldsymbol{\theta}) = (-y \log \hat{y} - (1-y) \log(1-\hat{y}))$$

- 4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \mathbf{x}(\hat{y} - y)$$

- 5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L \quad \eta \text{ is learning rate}$$

- 1) Pick a sample from training data
- 2) Compute output \hat{y}

$$z = \boldsymbol{\theta}^T \mathbf{x}$$

$$\mathbf{d} = [1 \dots 1] e^z$$

$$\hat{y} = e^z \emptyset \mathbf{d}$$

\emptyset is
Hadamard
division

- 3) Compute loss (cross-entropy)

$$L(\boldsymbol{\theta}) = -\mathbf{y}^T \log \hat{\mathbf{y}}$$

- 4) Compute derivative

$$\nabla_{\boldsymbol{\theta}} L = \mathbf{x}(\hat{\mathbf{y}} - \mathbf{y})^T$$

- 5) Update parameters

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \eta \nabla_{\boldsymbol{\theta}} L \quad \eta \text{ is learning rate}$$

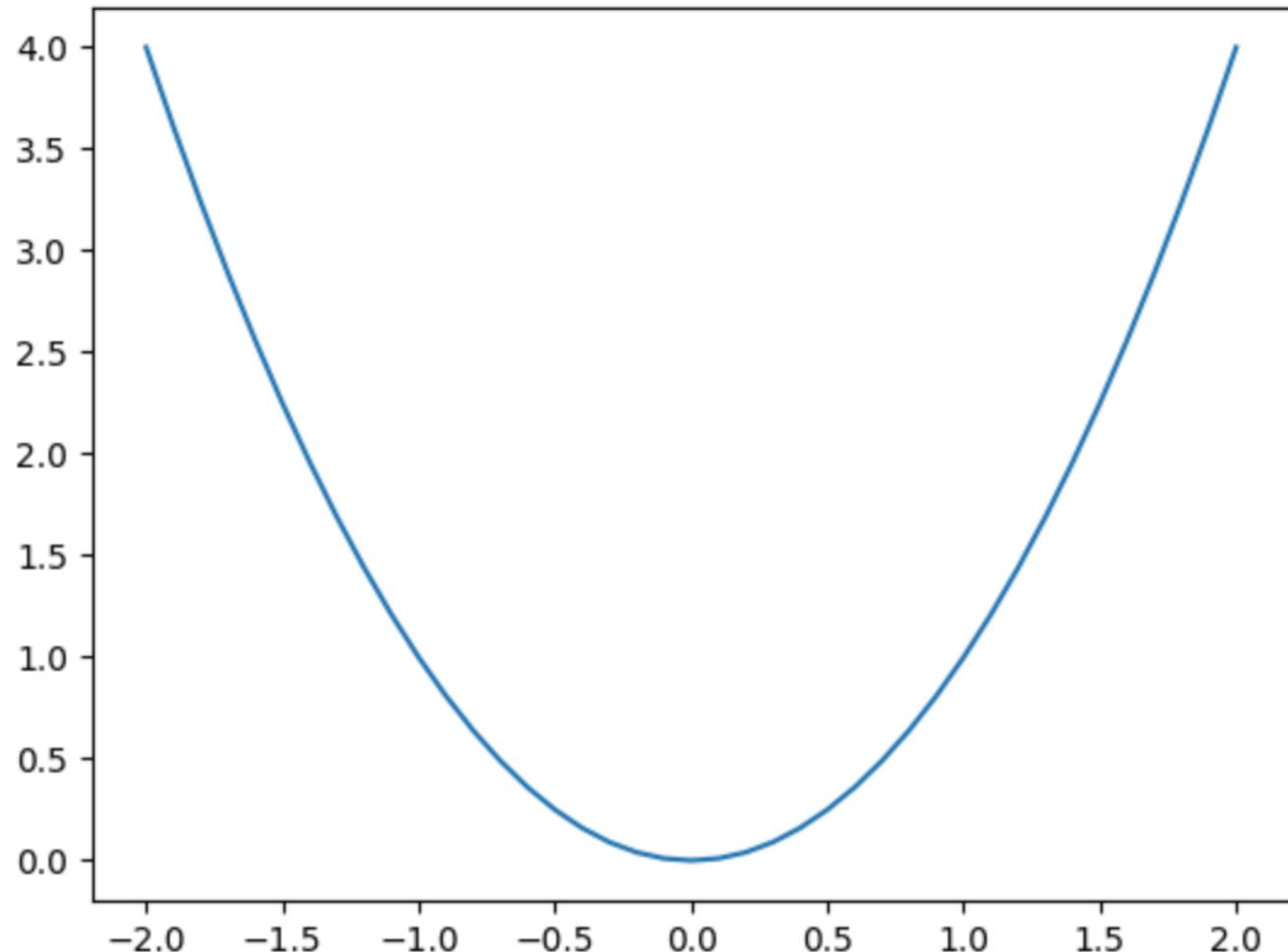
2 – GD with Momentum



Gradient Descent

- Objective function: $2x^2$
- Derivative: $2x$
- Gradient:

$$x = x - \eta * f'(x)$$



2 – GD with Momentum



Gradient Descent

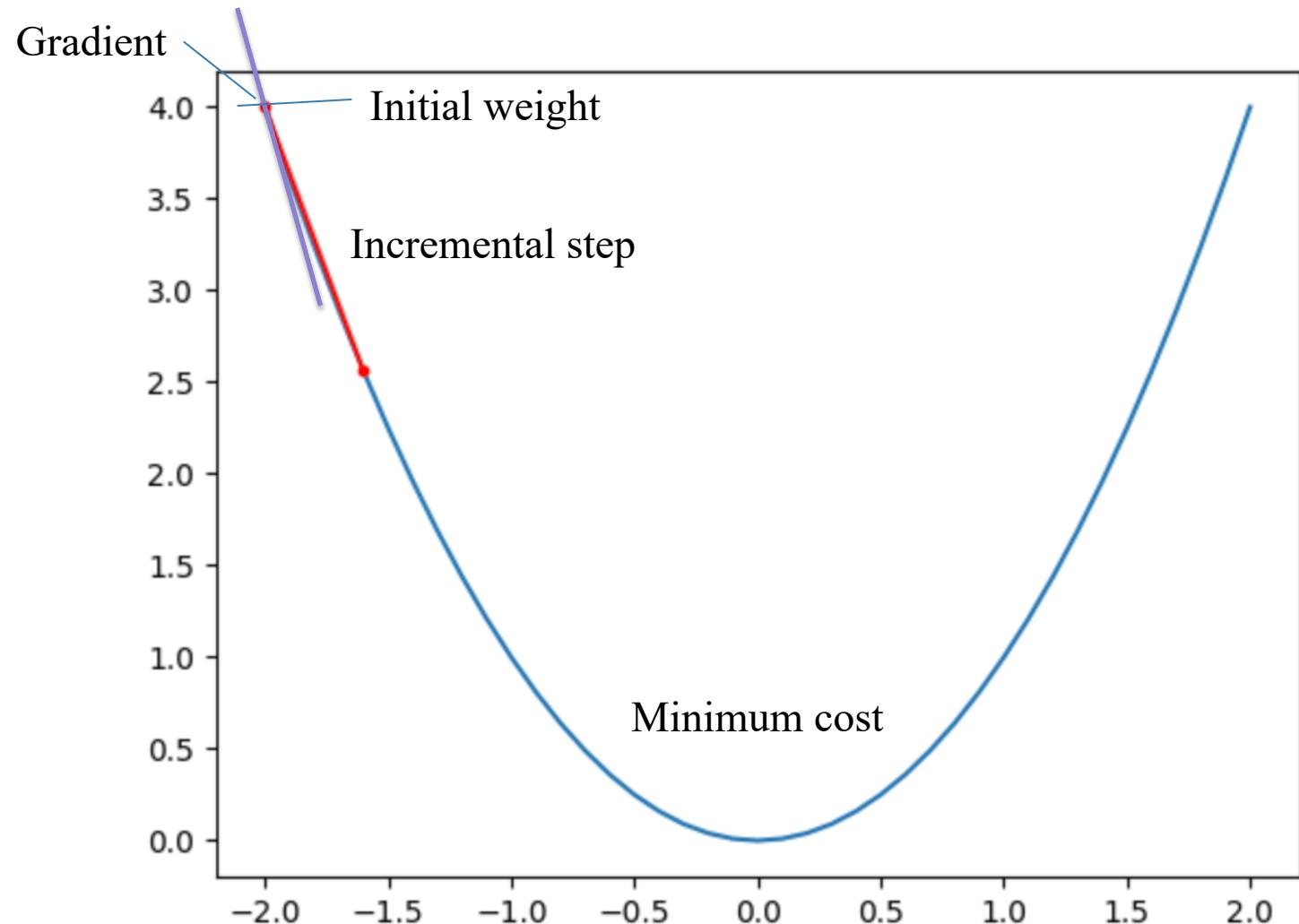
- Objective function: $2x^2$
- Derivative: $2x$
- Gradient:

$$x = x - \eta * f'(x)$$

$$x = -2.0$$

$$\eta = 0.1$$

$$\begin{aligned}x &= -2.0 - 0.1 * (2 * -2.0) \\&= -1.6\end{aligned}$$



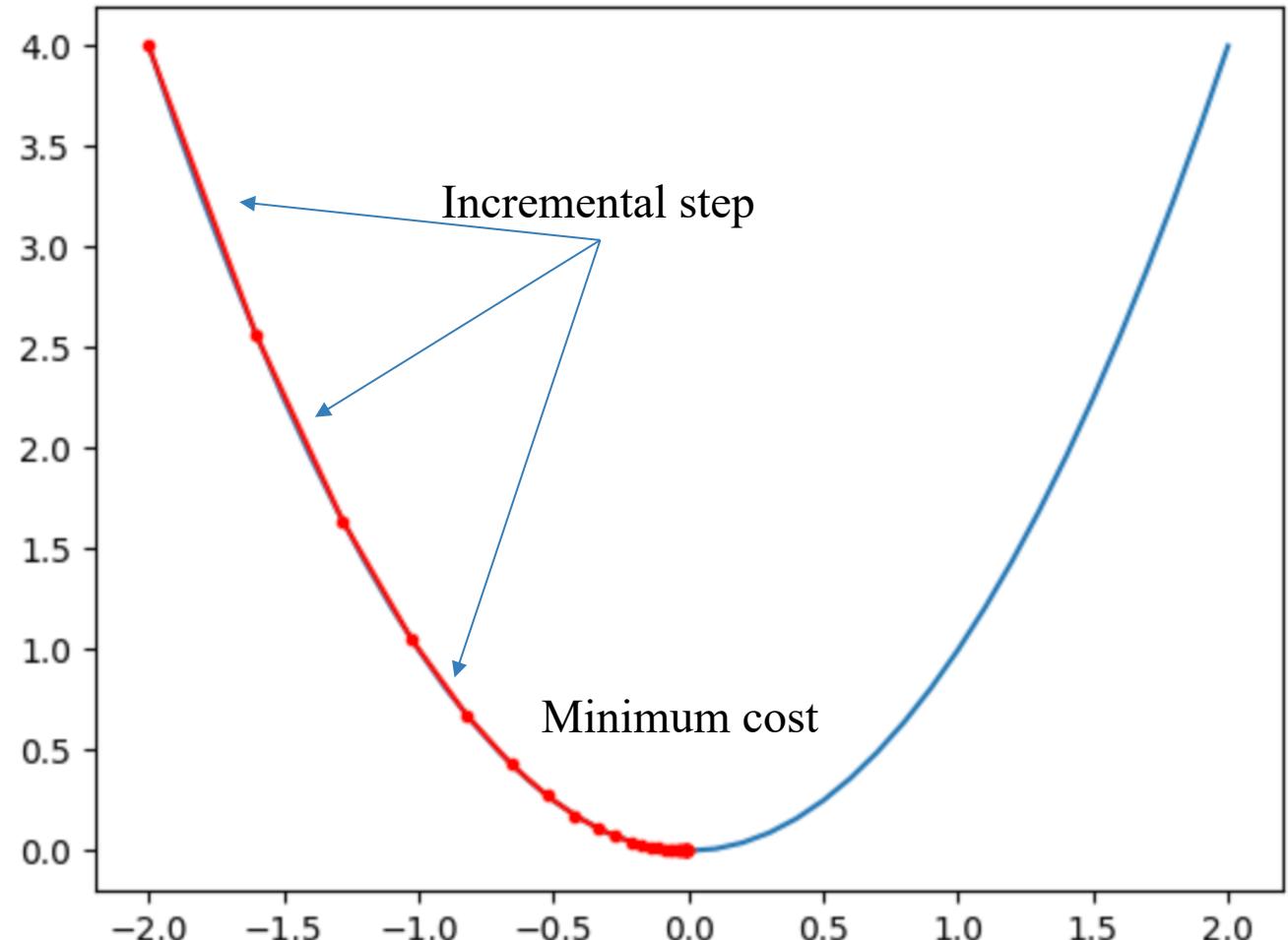
2 – GD with Momentum



Gradient Descent

- Objective function: $2x^2$
- Derivative: $2x$
- Gradient:

$$x = x - \eta * f'(x)$$



2 – GD with Momentum



Gradient Descent

```
# objective function
def objective(x):
    return x**2.0

# derivative of objective function
def derivative(x):
    return 2.0*x

# initial inputs
def init_inputs(r_min=-2.0, r_max=2.0):
    # sample input range uniformly at 0.1 increments
    inputs = arange(r_min, r_max+0.1, 0.1)
    return inputs
```

```
# gradient descent algorithm
def gradient_descent(inputs, num_epochs, learning_rate):
    # track all solutions
    solutions, scores = list(), list()
    # generate an initial point
    solution = inputs[0]
    # run the gradient descent
    solution_eval = objective(solution)
    # store solution
    solutions.append(solution)
    scores.append(solution_eval)
    for i in range(num_epochs):
        # calculate gradient
        gradient = derivative(solution)
        # take a step
        solution = solution - learning_rate * gradient
        # evaluate candidate point
        solution_eval = objective(solution)
        # store solution
        solutions.append(solution)
        scores.append(solution_eval)
        # report progress
        print('>%0.2d f(%0.5s) = %.5f' % (i, solution, solution_eval))
    return [solutions, scores]

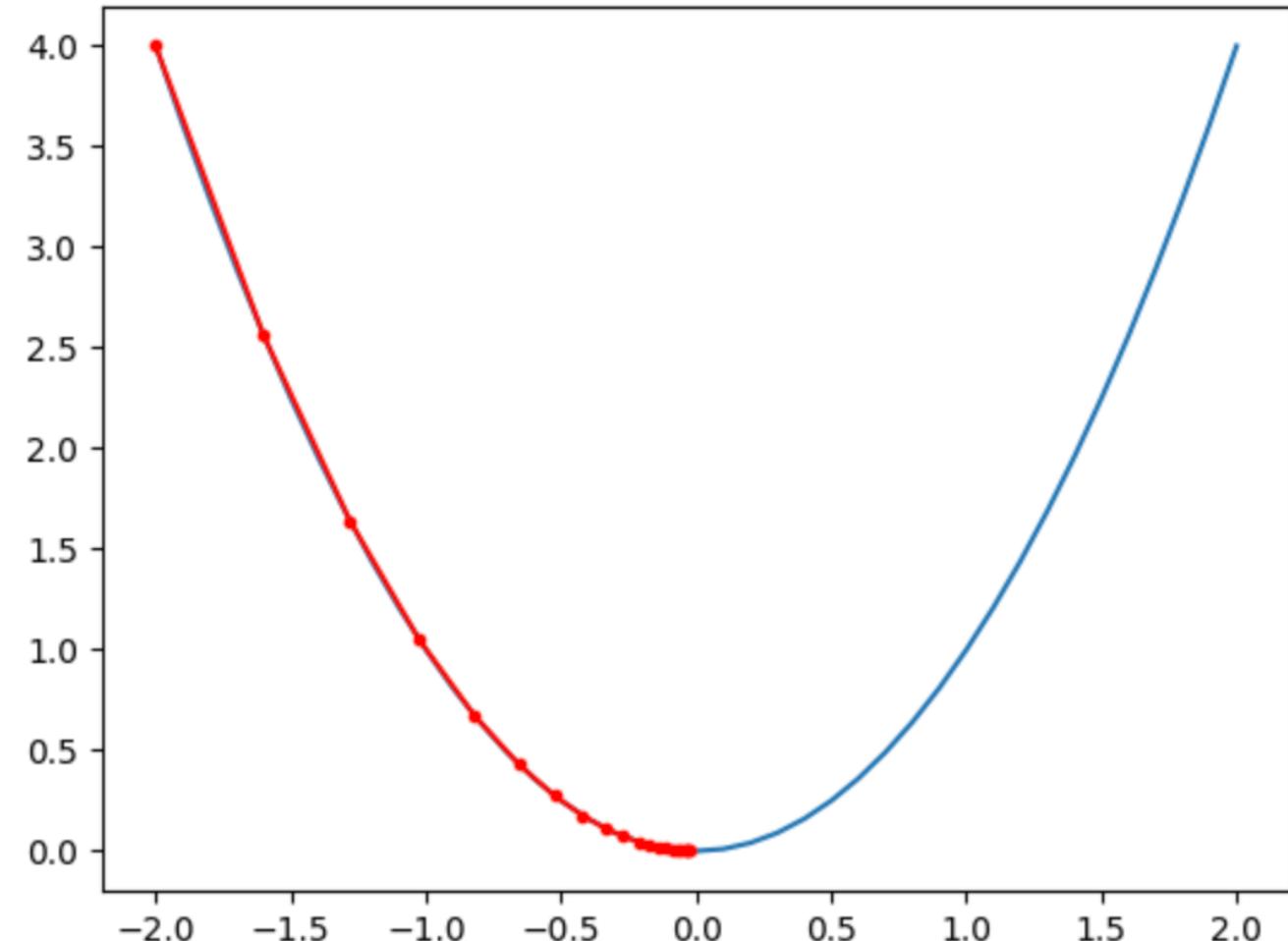
num_epochs = 20
learning_rate = 0.1
inputs = init_inputs()
solutions, scores = gradient_descent(inputs, num_epochs, learning_rate)
```

2 – GD with Momentum



Gradient Descent

```
Epoch: 00 -- f(-1.600) = 2.56000
Epoch: 01 -- f(-1.280) = 1.63840
Epoch: 02 -- f(-1.024) = 1.04858
Epoch: 03 -- f(-0.819) = 0.67109
Epoch: 04 -- f(-0.655) = 0.42950
Epoch: 05 -- f(-0.524) = 0.27488
Epoch: 06 -- f(-0.419) = 0.17592
Epoch: 07 -- f(-0.336) = 0.11259
Epoch: 08 -- f(-0.268) = 0.07206
Epoch: 09 -- f(-0.215) = 0.04612
Epoch: 10 -- f(-0.172) = 0.02951
Epoch: 11 -- f(-0.137) = 0.01889
Epoch: 12 -- f(-0.110) = 0.01209
Epoch: 13 -- f(-0.088) = 0.00774
Epoch: 14 -- f(-0.070) = 0.00495
Epoch: 15 -- f(-0.056) = 0.00317
Epoch: 16 -- f(-0.045) = 0.00203
Epoch: 17 -- f(-0.036) = 0.00130
Epoch: 18 -- f(-0.029) = 0.00083
Epoch: 19 -- f(-0.023) = 0.00053
```

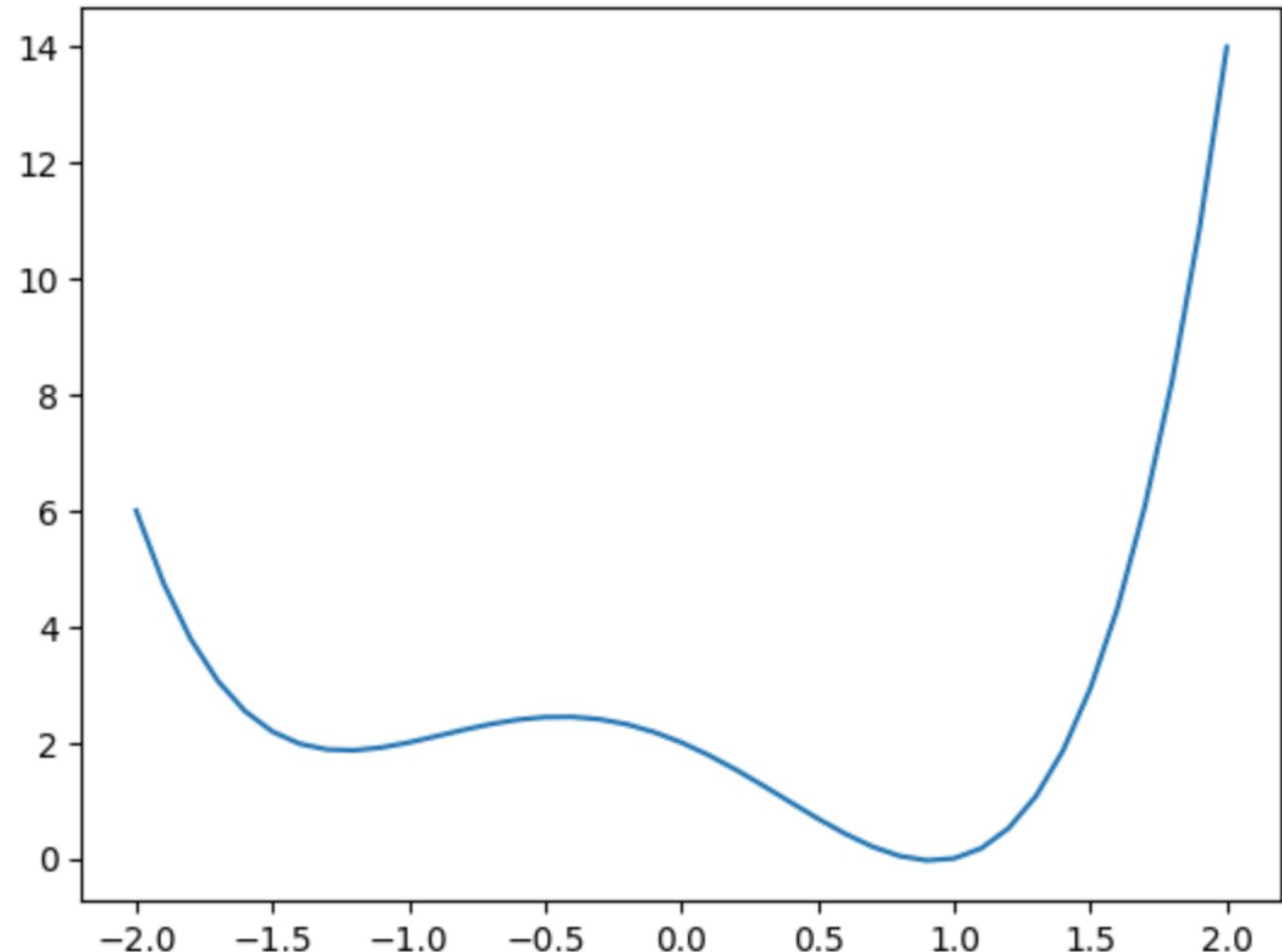


2 – GD with Momentum



Gradient Descent

- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient:
 $x = x - \eta * f'(x)$



2 – GD with Momentum



Gradient Descent

- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$

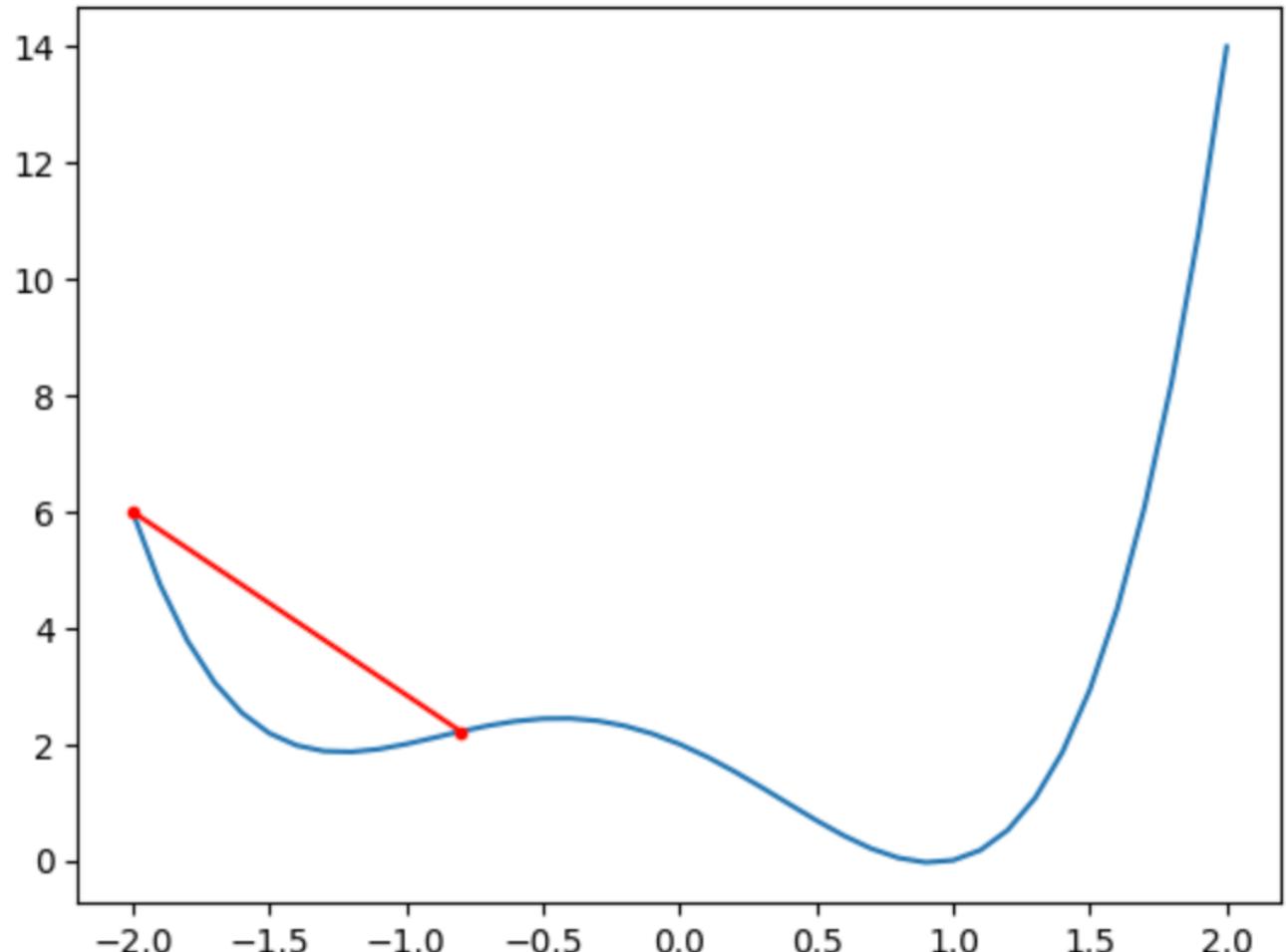
- Derivative:
 $4x^3 + 3x^2 - 4x$

- Gradient:
$$x = x - \eta * f'(x)$$

$$x = -2.0$$

$$\eta = 0.1$$

$$\begin{aligned}x &= -2.0 - 0.1 * (-12) \\&= -0.8\end{aligned}$$



2 – GD with Momentum



Gradient Descent

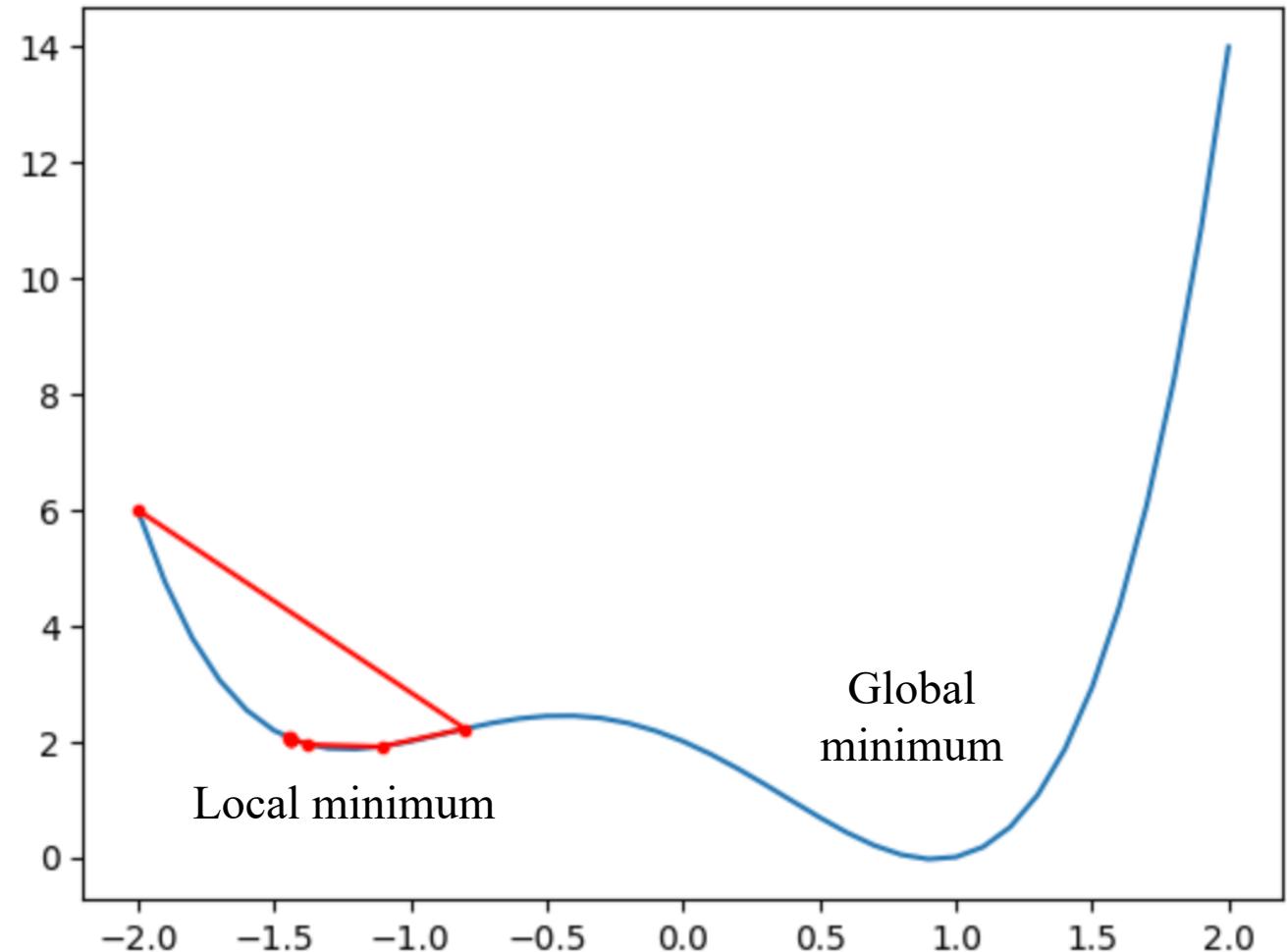
- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient:
 $x = x - \eta * f'(x)$

$x = -2.0$

$\eta = 0.1$

num_epochs = 5

Cost = 2.05



2 – GD with Momentum



Gradient Descent

```
# objective function
def objective(x):
    return x**4+x**3-2*x**2-2*x + 2

# derivative of objective function
def derivative(x):
    return 4*x**3 + 3*x**2 - 4*x

# initial inputs
def init_inputs(r_min=-2.0, r_max=2.0):
    # sample input range uniformly at 0.1 increments
    inputs = arange(r_min, r_max+0.1, 0.1)
    return inputs
```

```
# gradient descent algorithm
def gradient_descent(inputs, num_epochs, learning_rate):
    # track all solutions
    solutions, scores = list(), list()
    # generate an initial point
    solution = inputs[0]
    # run the gradient descent
    solution_eval = objective(solution)
    # store solution
    solutions.append(solution)
    scores.append(solution_eval)
    # run the gradient descent
    for i in range(num_epochs):
        # calculate gradient
        gradient = derivative(solution)
        # take a step
        solution = solution - learning_rate * gradient
        # evaluate candidate point
        solution_eval = objective(solution)
        # store solution
        solutions.append(solution)
        scores.append(solution_eval)
        # report progress
        print('Epoch: %0.2d -- f(%0.3f) = %.5f' % (i, solution, solution_eval))
    return [solutions, scores]

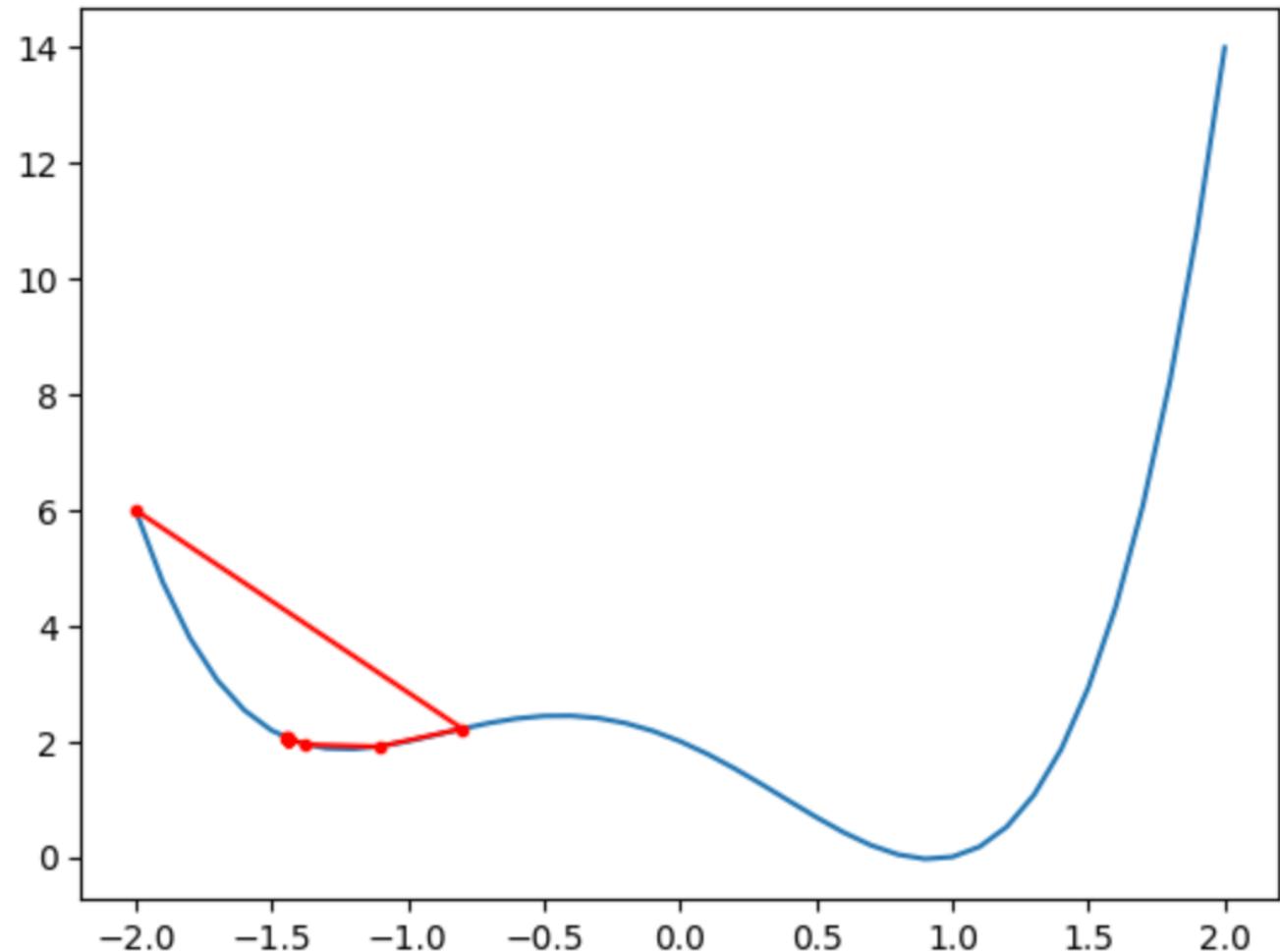
num_epochs = 20
learning_rate = 0.1
inputs = init_inputs()
solutions, scores = gradient_descent(inputs, num_epochs, learning_rate)
```

2 – GD with Momentum



Gradient Descent

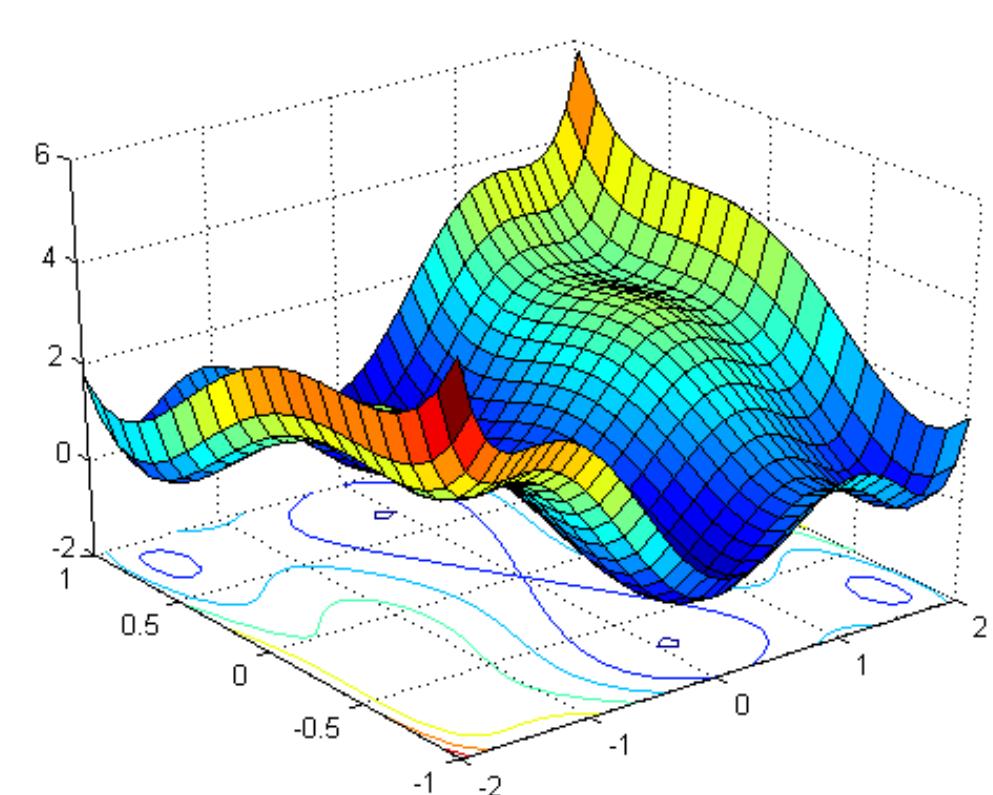
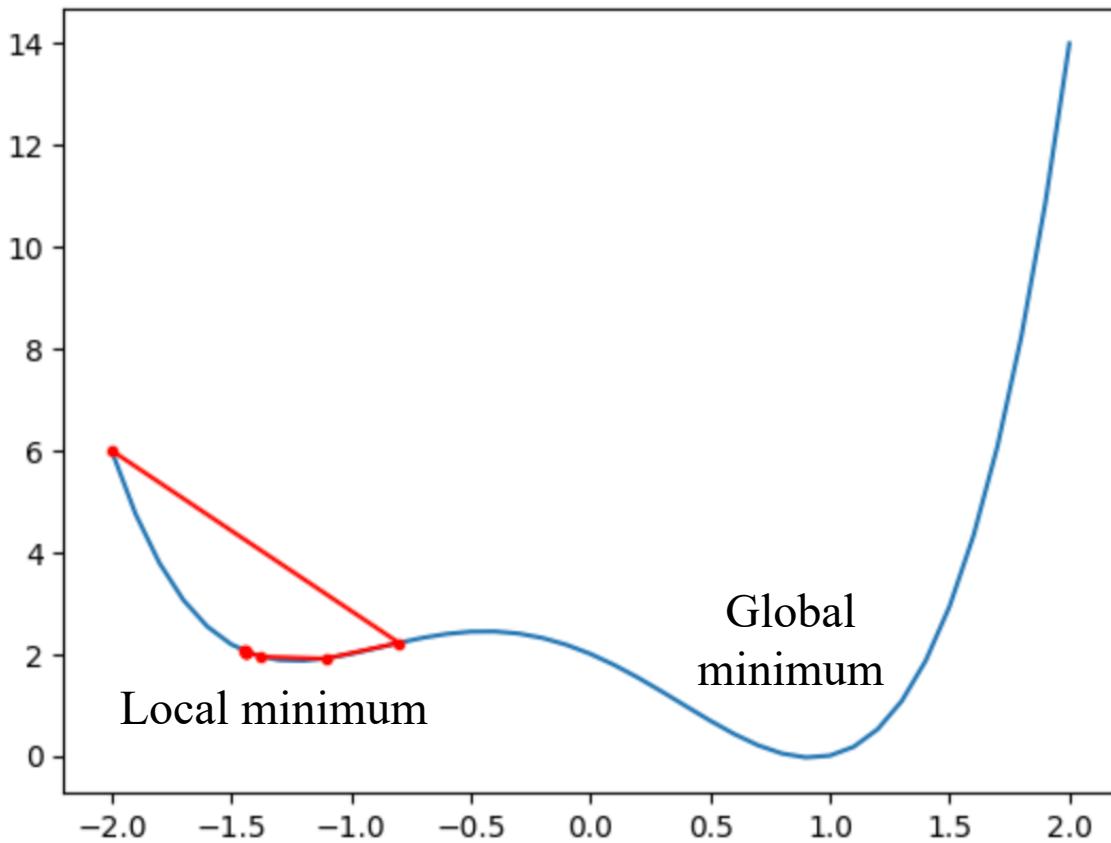
```
Epoch: 00 -- f(-0.800) = 2.21760
Epoch: 01 -- f(-1.107) = 1.90812
Epoch: 02 -- f(-1.375) = 1.94351
Epoch: 03 -- f(-1.452) = 2.07181
Epoch: 04 -- f(-1.441) = 2.04801
Epoch: 05 -- f(-1.444) = 2.05364
Epoch: 06 -- f(-1.443) = 2.05233
Epoch: 07 -- f(-1.443) = 2.05264
Epoch: 08 -- f(-1.443) = 2.05257
Epoch: 09 -- f(-1.443) = 2.05258
Epoch: 10 -- f(-1.443) = 2.05258
Epoch: 11 -- f(-1.443) = 2.05258
Epoch: 12 -- f(-1.443) = 2.05258
Epoch: 13 -- f(-1.443) = 2.05258
Epoch: 14 -- f(-1.443) = 2.05258
Epoch: 15 -- f(-1.443) = 2.05258
Epoch: 16 -- f(-1.443) = 2.05258
Epoch: 17 -- f(-1.443) = 2.05258
Epoch: 18 -- f(-1.443) = 2.05258
Epoch: 19 -- f(-1.443) = 2.05258
```



2 – GD with Momentum



Gradient Descent



2 – GD with Momentum

!

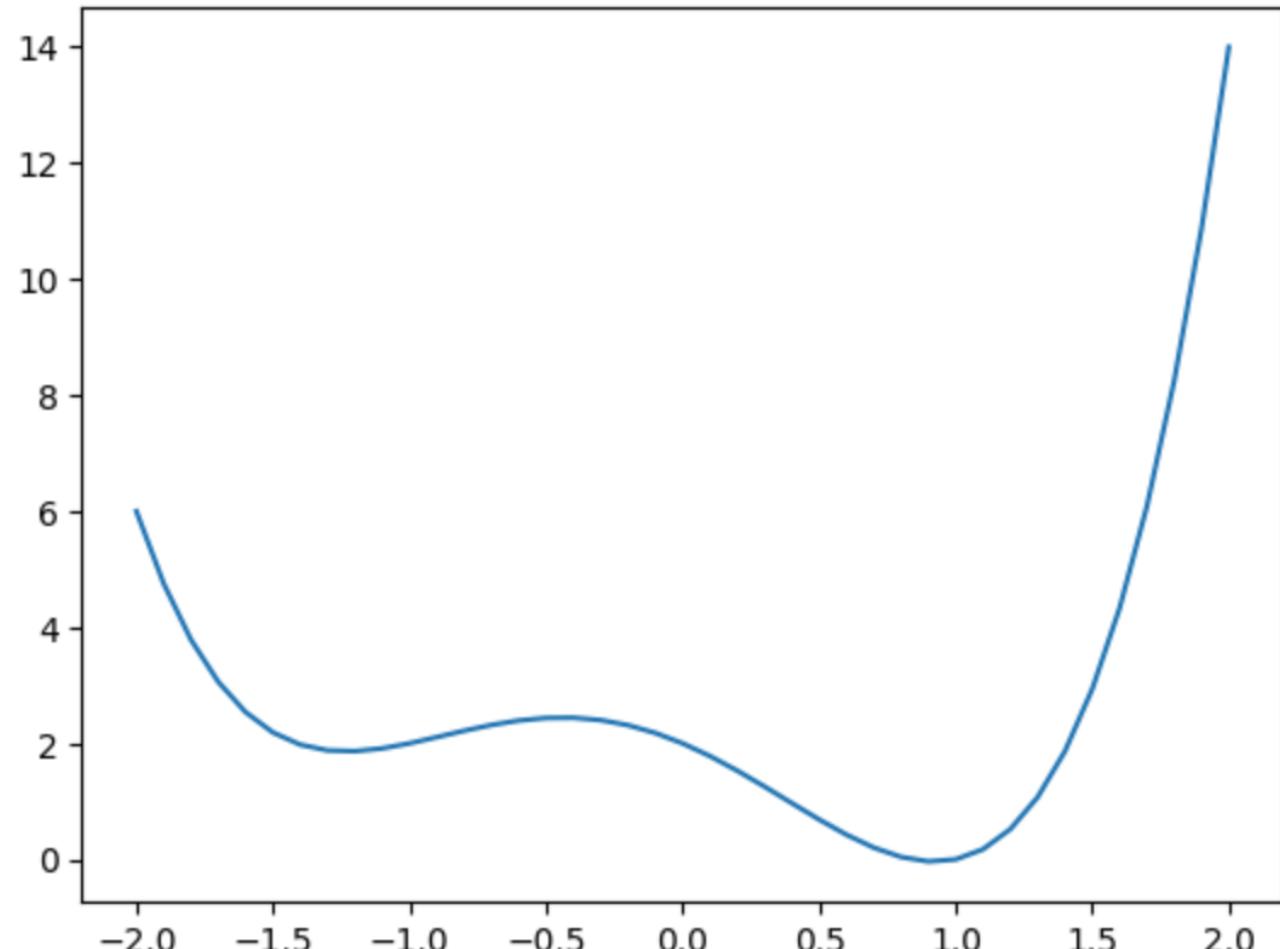
GD with Momentum

- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient with Momentum:

$$v_t = \gamma v_{t-1} + \eta f'(x)$$

$$x = x - v_t$$

γ : momentum



2 – GD with Momentum



GD with Momentum

- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient with Momentum:

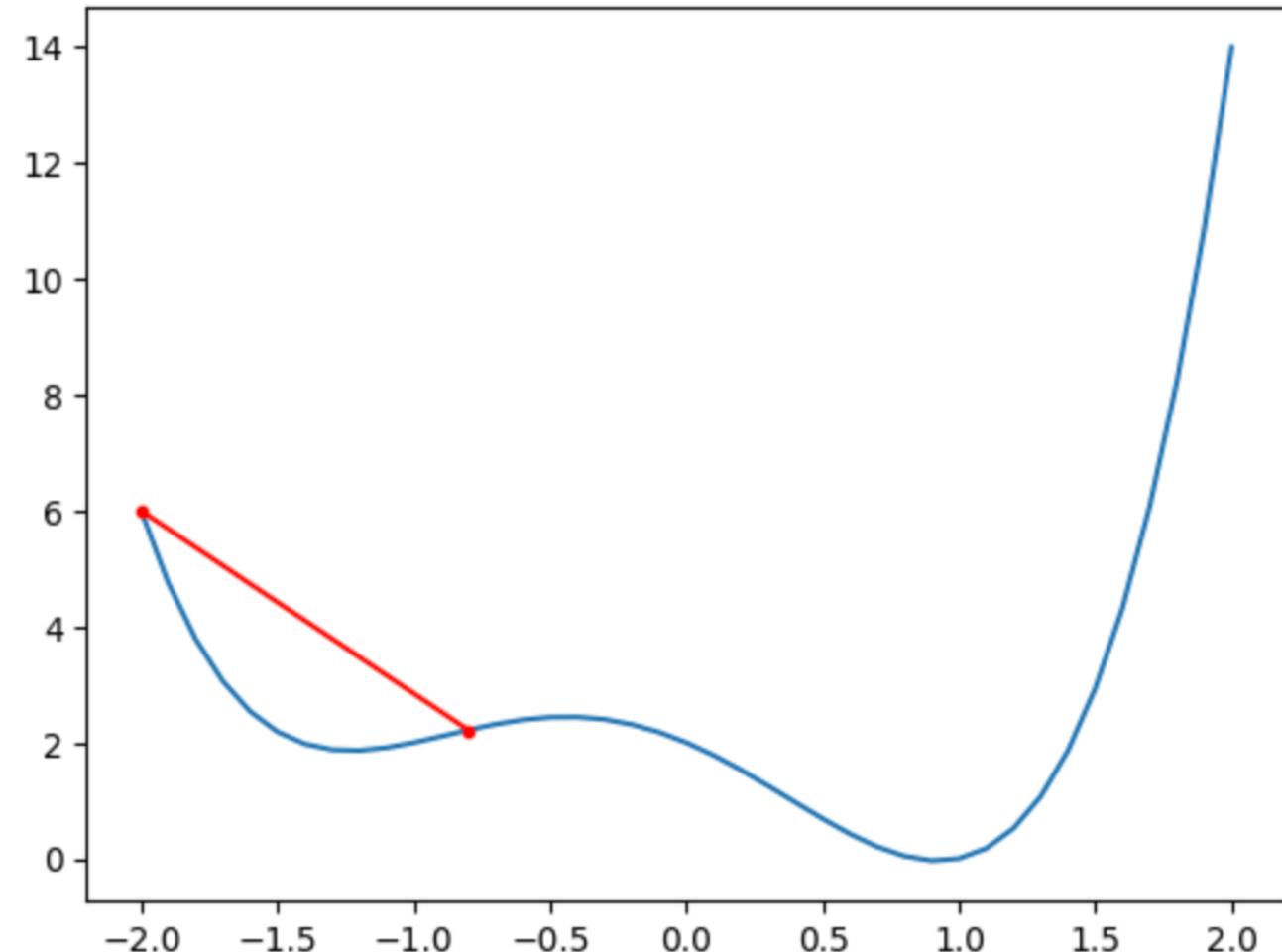
$$v_t = \gamma v_{t-1} + \eta f'(x)$$

$$x = x - v_t$$

$$x = -2.0 \quad \eta = 0.1 \quad \gamma = 0.8 \quad v_0 = 0.0$$

$$v_1 = 0.8 * 0.0 + 0.1 * (-12) = -1.2$$

$$x = -2.0 - (-1.2) = -0.8$$



2 – GD with Momentum



GD with Momentum

- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient with Momentum:

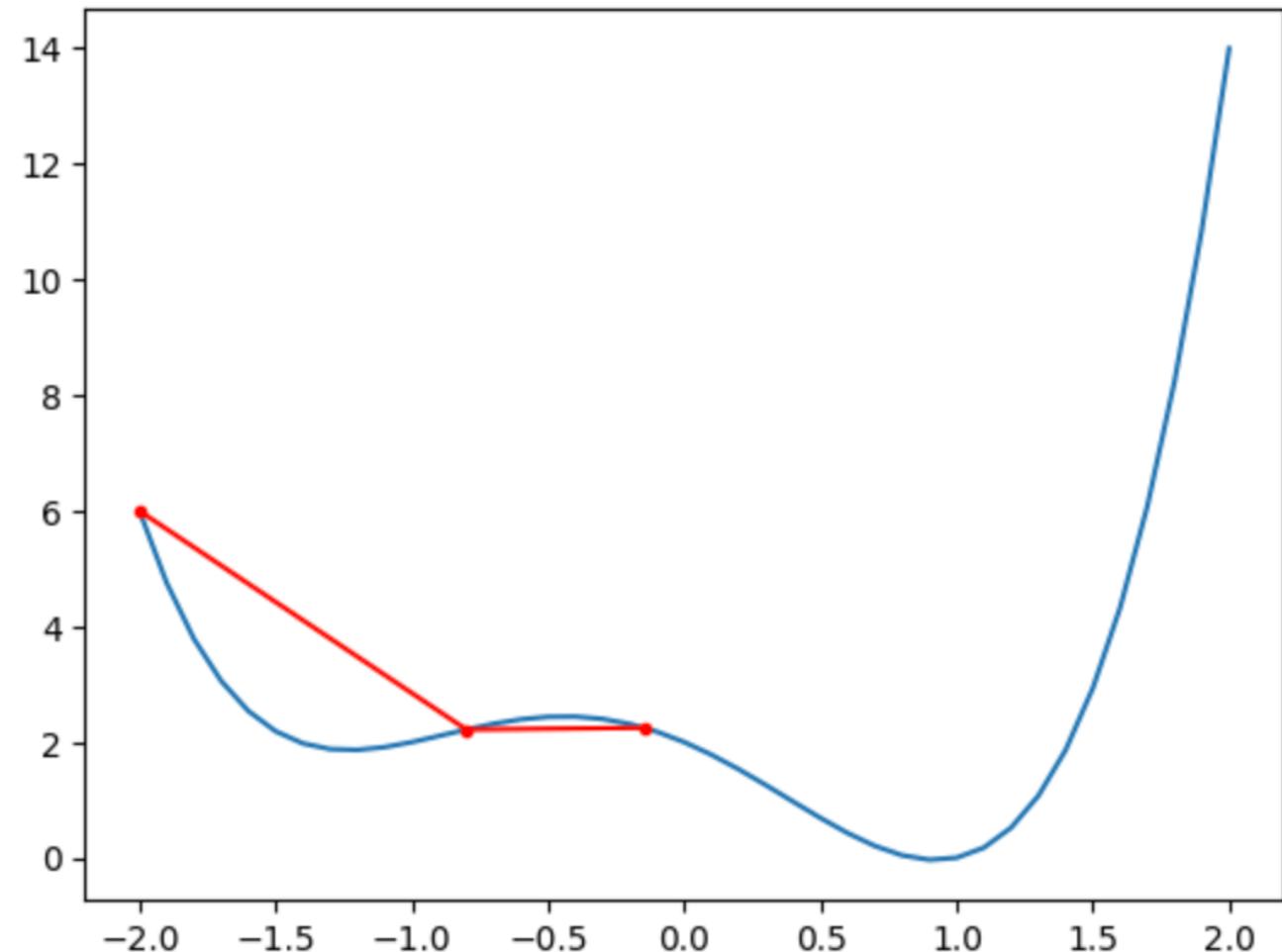
$$v_t = \gamma v_{t-1} + \eta f'(x)$$

$$x = x - v_t$$

$$x = -2.0 \quad \eta = 0.1 \quad \gamma = 0.8 \quad v_1 = -1.2$$

$$v_1 = 0.8 * (-1.2) + 0.1 * (3.072) = -0.65$$

$$x = -0.8 - (-0.65) = -0.15$$



2 – GD with Momentum

!

GD with Momentum

- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient with Momentum:

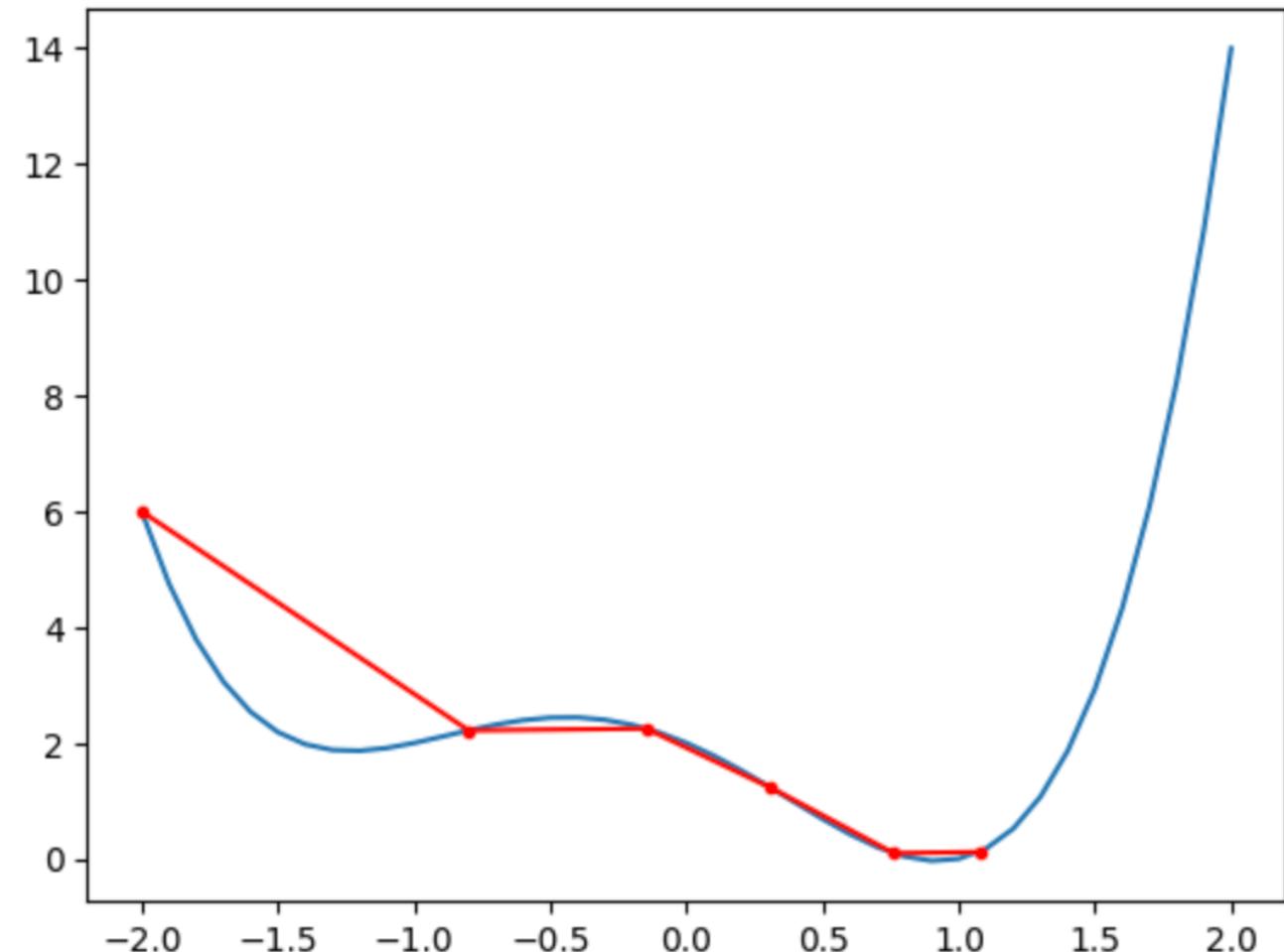
$$v_t = \gamma v_{t-1} + \eta f'(x)$$

$$x = x - v_t$$

$$x = -2.0 \quad \eta = 0.1 \quad \gamma = 0.8 \quad v_0 = 0.0$$

$$\text{num_epochs} = 5$$

$$\text{Cost} = 0.12$$



2 – GD with Momentum

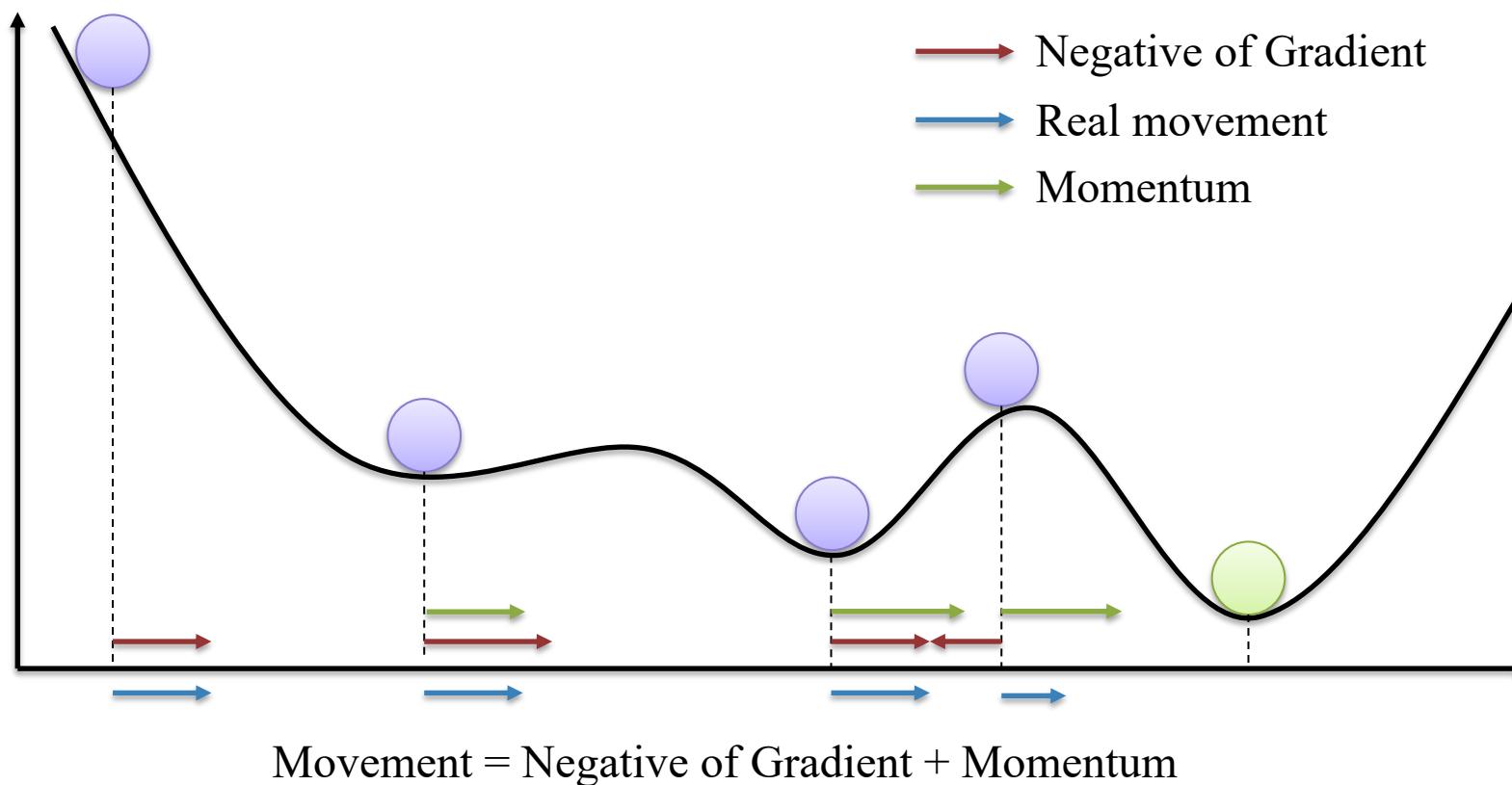


GD with Momentum

- Gradient with Momentum:

$$v_t = \gamma v_{t-1} + \eta f'(x)$$

$$x = x - v_t$$



2 – GD with Momentum



GD with Momentum

```
# objective function
def objective(x):
    return x**4+x**3-2*x**2-2*x + 2

# derivative of objective function
def derivative(x):
    return 4*x**3 + 3*x**2 - 4*x

# initial inputs
def init_inputs(r_min=-2.0, r_max=2.0):
    # sample input range uniformly at 0.1 increments
    inputs = arange(r_min, r_max+0.1, 0.1)
    return inputs

num_epochs = 20
learning_rate = 0.1
inputs = init_inputs()
# define momentum
momentum = 0.8
solutions, scores = gradient_descent_with_momentum(
    inputs, num_epochs, learning_rate, momentum
)
```

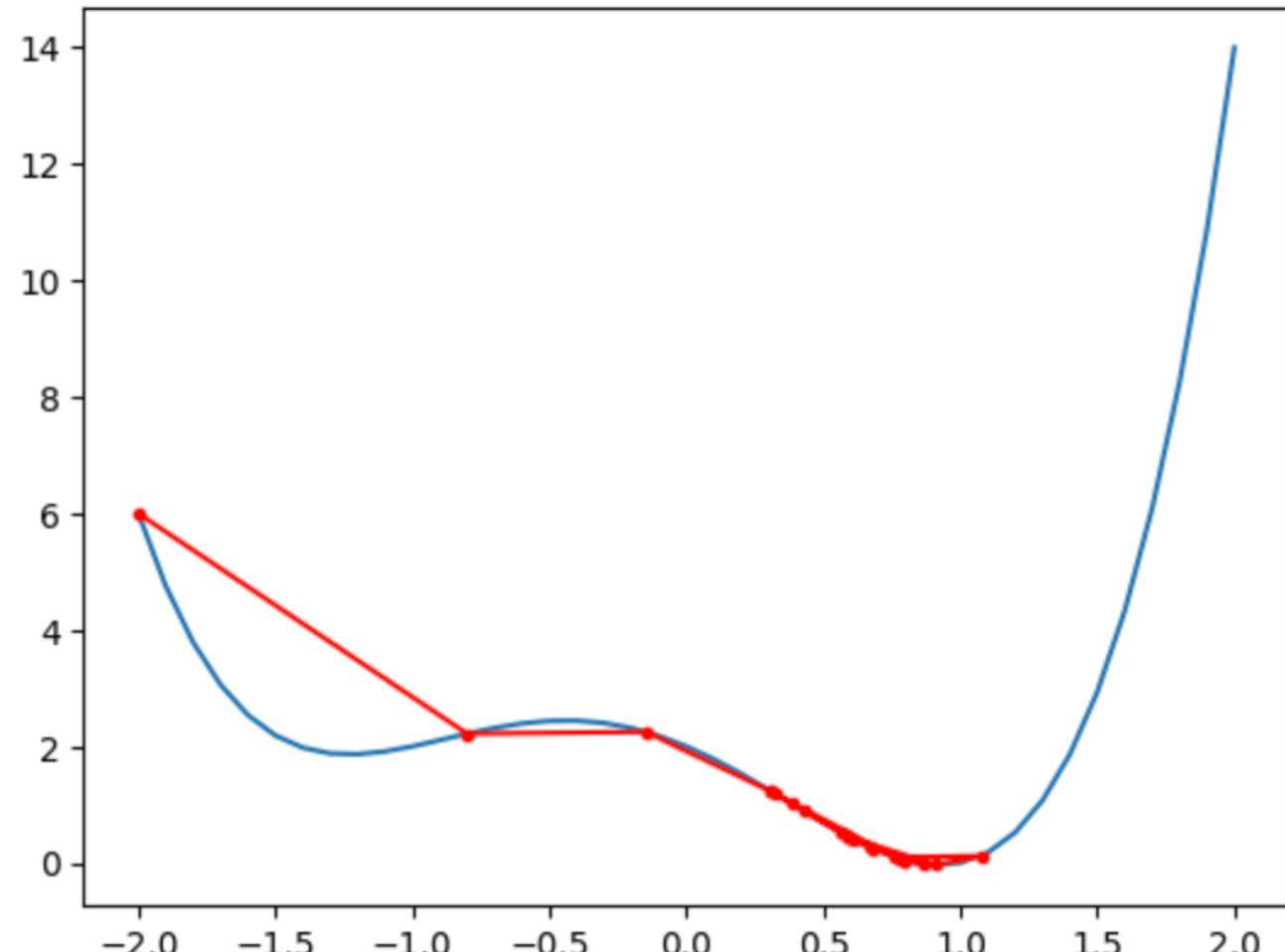
```
# gradient descent algorithm
def gradient_descent_with_momentum(inputs, num_epochs, learning_rate, momentum):
    # track all solutions
    solutions, scores = list(), list()
    # generate an initial point
    solution = inputs[0]
    # keep track of the change
    change = 0.0
    # run the gradient descent
    solution_eval = objective(solution)
    # store solution
    solutions.append(solution)
    scores.append(solution_eval)
    for i in range(num_epochs):
        # calculate gradient
        gradient = derivative(solution)
        # calculate update
        new_change = learning_rate * gradient + momentum * change
        # take a step
        solution = solution - new_change
        # save the change
        change = new_change
        # evaluate candidate point
        solution_eval = objective(solution)
        # store solution
        solutions.append(solution)
        scores.append(solution_eval)
        # report progress
        print('Epoch: %0.2d -- f(%0.3f) = %.5f' % (i, solution, solution_eval))
    return [solutions, scores]
```

2 – GD with Momentum

!

GD with Momentum

```
Epoch: 00 -- f(-0.800) = 2.21760
Epoch: 01 -- f(-0.147) = 2.24834
Epoch: 02 -- f(0.311) = 1.22418
Epoch: 03 -- f(0.761) = 0.09618
Epoch: 04 -- f(1.075) = 0.11696
Epoch: 05 -- f(0.913) = -0.03723
Epoch: 06 -- f(0.594) = 0.44088
Epoch: 07 -- f(0.387) = 1.00798
Epoch: 08 -- f(0.308) = 1.23381
Epoch: 09 -- f(0.327) = 1.17773
Epoch: 10 -- f(0.428) = 0.89005
Epoch: 11 -- f(0.593) = 0.44244
Epoch: 12 -- f(0.774) = 0.07689
Epoch: 13 -- f(0.863) = -0.01795
Epoch: 14 -- f(0.799) = 0.04293
Epoch: 15 -- f(0.672) = 0.26025
Epoch: 16 -- f(0.582) = 0.46945
Epoch: 17 -- f(0.563) = 0.51926
Epoch: 18 -- f(0.606) = 0.41078
Epoch: 19 -- f(0.684) = 0.23563
```



3 – Nesterov Momentum



Problem of GD with Momentum

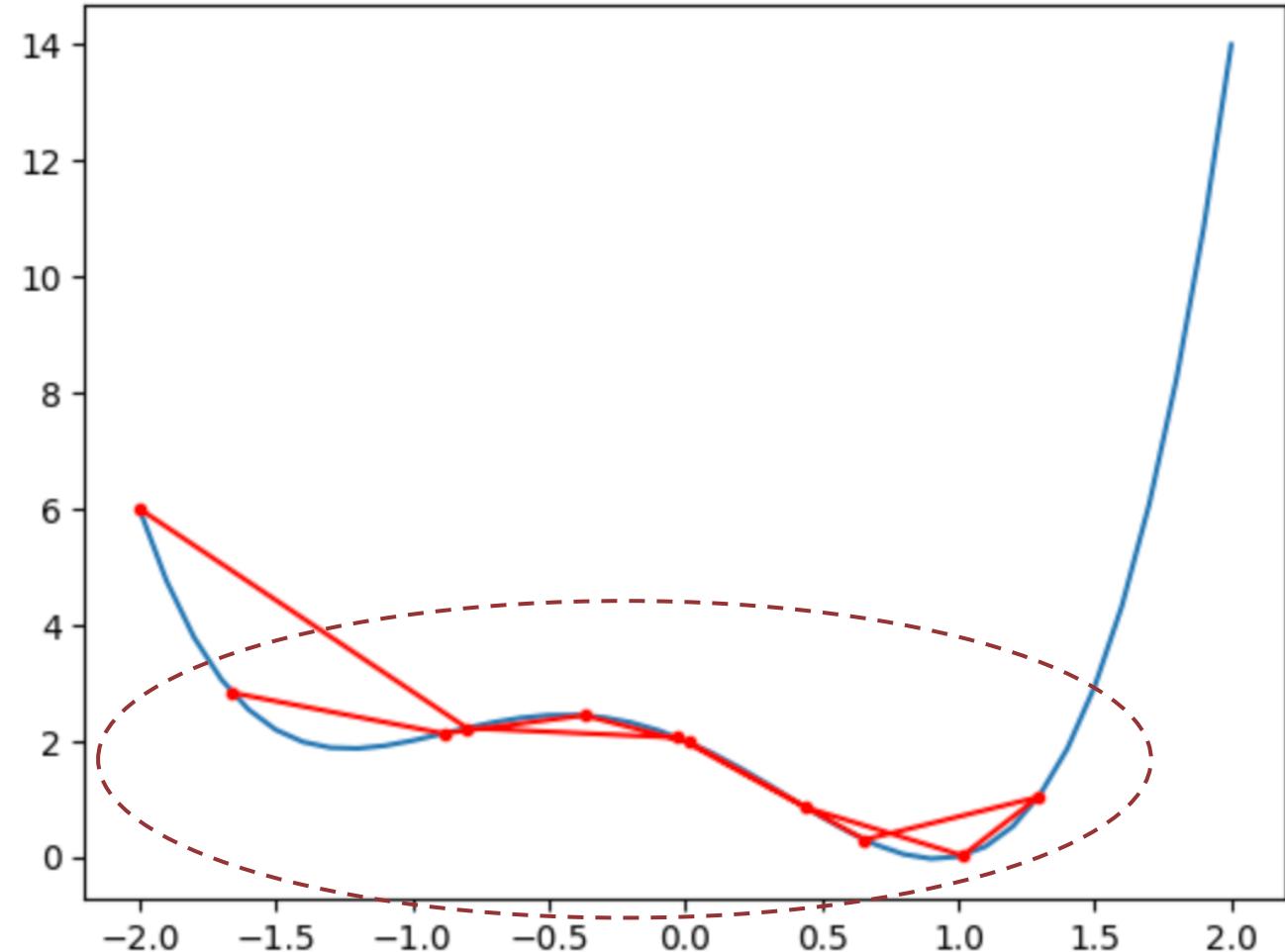
- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient with Momentum:

$$v_t = \gamma v_{t-1} + \eta f'(x)$$

$$x = x - v_t$$

$$x = -2.0 \quad \eta = 0.1 \quad \gamma = 0.9 \quad v_0 = 0.0$$

$$\text{num_epochs} = 10$$



3 – Nesterov Momentum

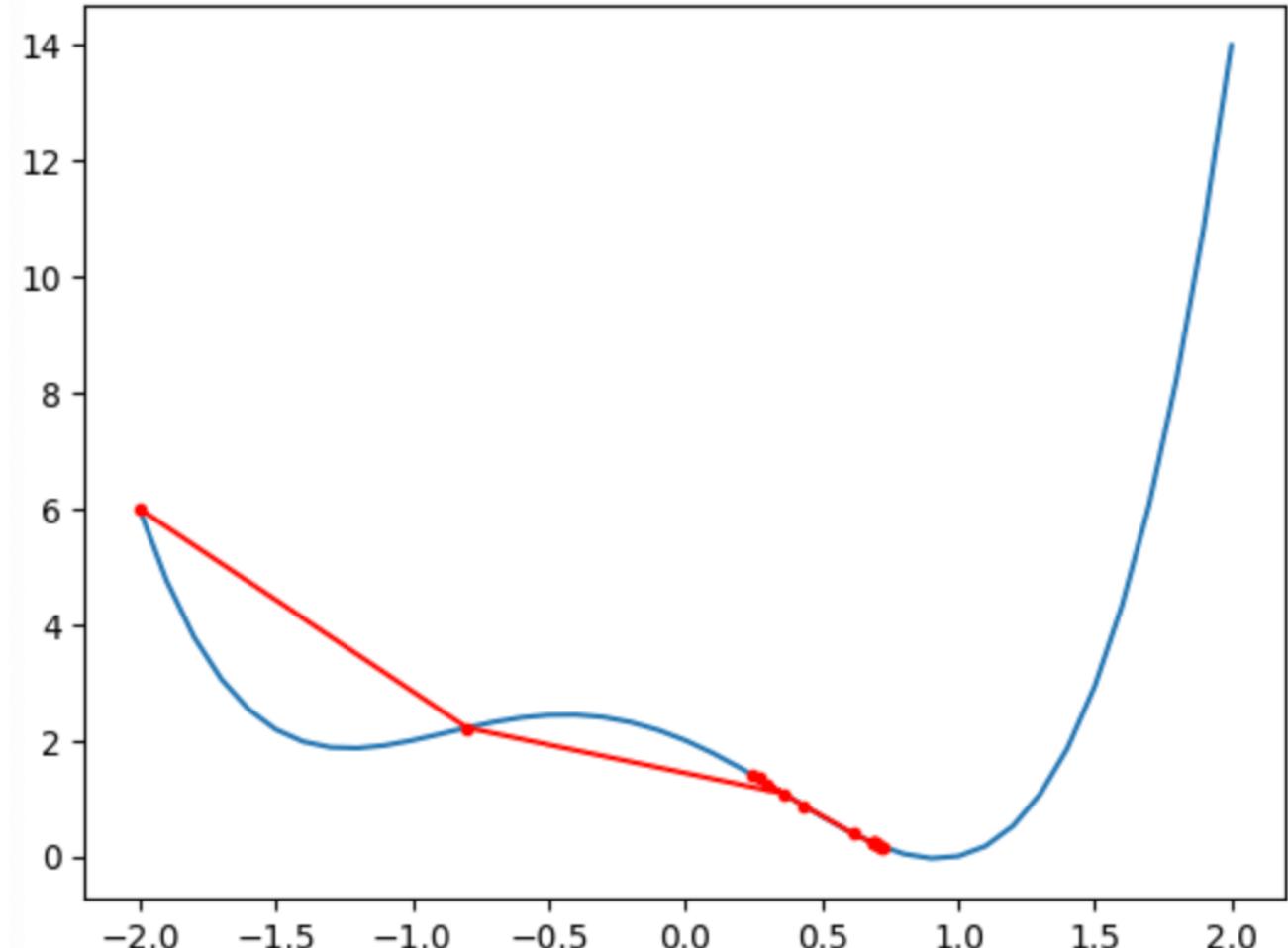


Nesterov Momentum

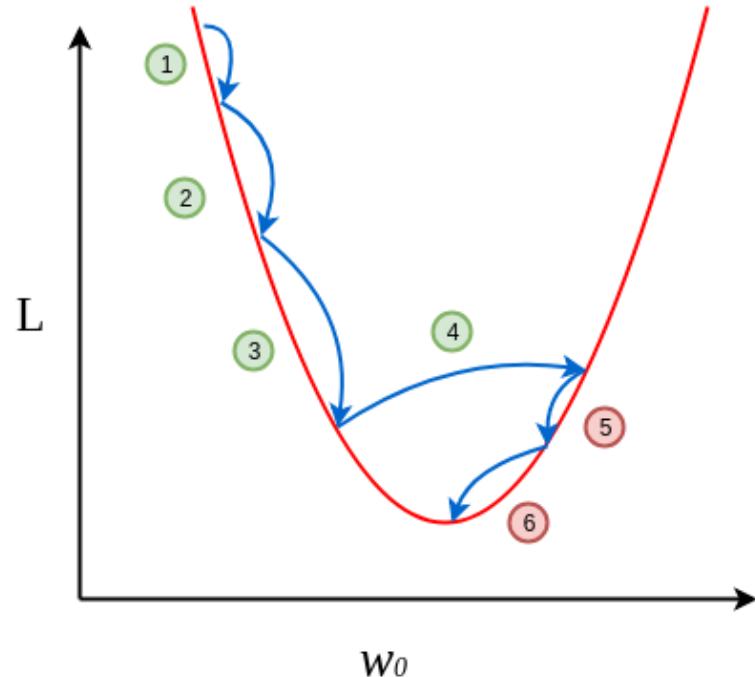
- Objective function:
 $x^4 + x^3 - 2x^2 - 2x + 2$
- Derivative:
 $4x^3 + 3x^2 - 4x$
- Gradient with Momentum:

$$v_t = \gamma v_{t-1} - \eta f'(x + \gamma * v_{t-1})$$
$$x = x + v_t$$

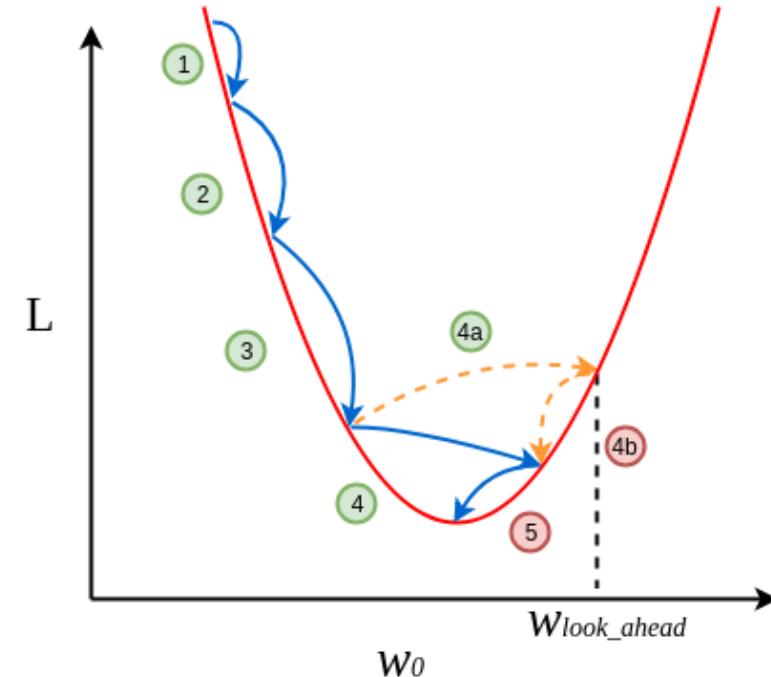
$$x = -2.0 \quad \eta = 0.1 \quad \gamma = 0.9 \quad v_0 = 0.0$$



3 – Nesterov Momentum



(a) Momentum-Based Gradient Descent



(b) Nesterov Accelerated Gradient Descent

Source

$$\text{Green Circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Positive}(+)}$$

$$\text{Red Circle} \Rightarrow \frac{\partial L}{\partial w_0} = \frac{\text{Negative}(-)}{\text{Negative}(-)}$$

3 – Nesterov Momentum



```
# objective function
def objective(x):
    return x**4+x**3-2*x**2-2*x + 2

# derivative of objective function
def derivative(x):
    return 4*x**3 + 3*x**2 - 4*x

# initial inputs
def init_inputs(r_min=-2.0, r_max=2.0):
    # sample input range uniformly at 0.1 increments
    inputs = arange(r_min, r_max+0.1, 0.1)
    return inputs

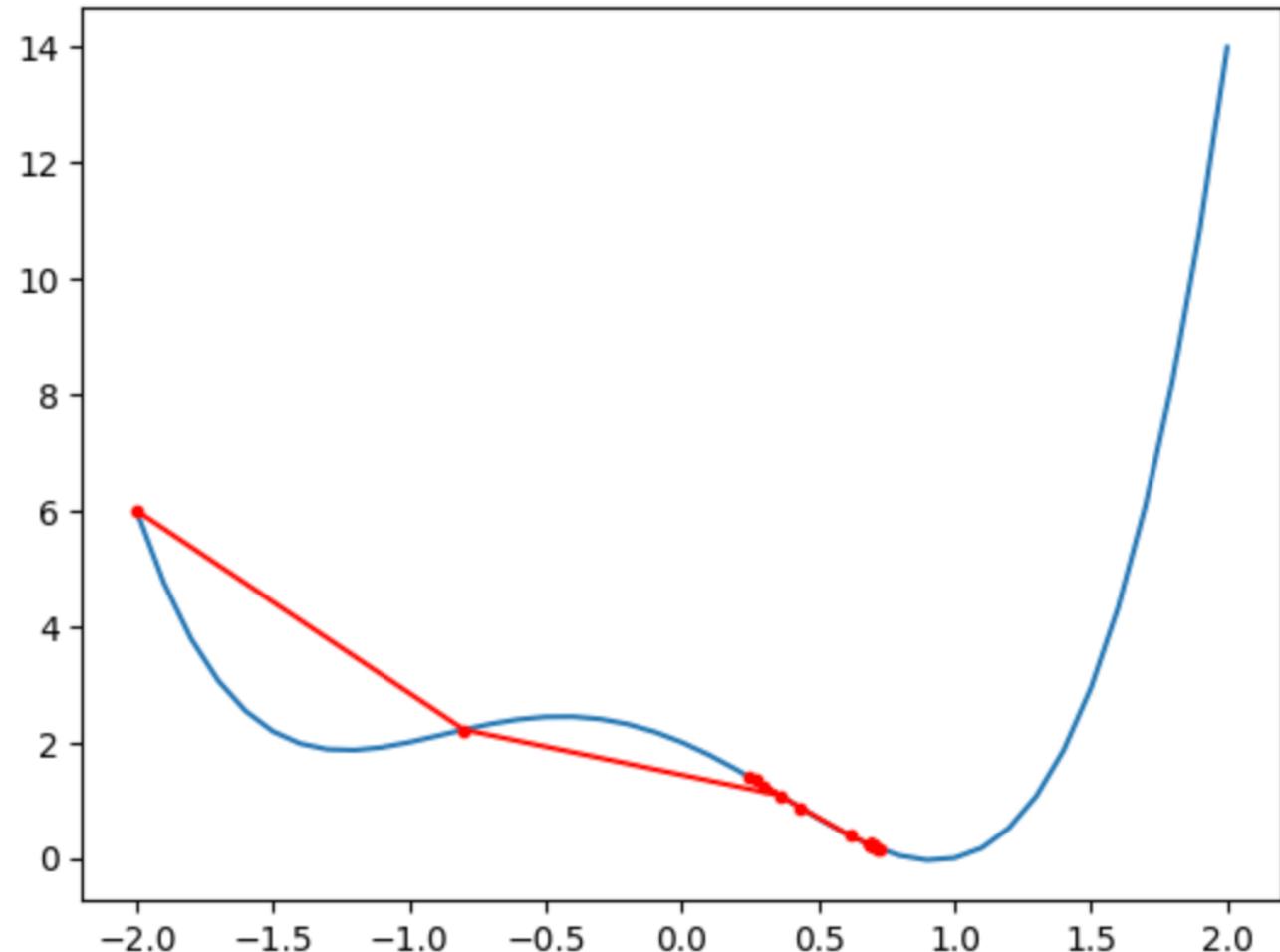
num_epochs = 20
learning_rate = 0.1
inputs = init_inputs()
# define momentum
momentum = 0.9
solutions, scores = nesterov_momentum(
    inputs, num_epochs, learning_rate, momentum
)
```

```
# gradient descent algorithm
def nesterov_momentum(inputs, num_epochs, learning_rate, momentum):
    # track all solutions
    solutions, scores = list(), list()
    # generate an initial point
    solution = inputs[0]
    # keep track of the change
    change = 0.0
    # run the gradient descent
    solution_eval = objective(solution)
    # store solution
    solutions.append(solution)
    scores.append(solution_eval)
    for i in range(num_epochs):
        # calculate the projected solution
        projected = solution + momentum * change
        # calculate gradient
        gradient = derivative(projected)
        # calculate update
        new_change = momentum * change - learning_rate * gradient
        # take a step
        solution = solution + new_change
        # save the change
        change = new_change
        # evaluate candidate point
        solution_eval = objective(solution)
        # store solution
        solutions.append(solution)
        scores.append(solution_eval)
        # report progress
        print('Epoch: %0.2d -- f(%0.3f) = %.5f' % (i, solution, solution_eval))
    return [solutions, scores]
```

3 – Nesterov Momentum

!

```
Epoch: 00 -- f(-0.800) = 2.21760
Epoch: 01 -- f(0.360) = 1.08511
Epoch: 02 -- f(0.268) = 1.34411
Epoch: 03 -- f(0.247) = 1.40184
Epoch: 04 -- f(0.300) = 1.25650
Epoch: 05 -- f(0.432) = 0.87708
Epoch: 06 -- f(0.614) = 0.39136
Epoch: 07 -- f(0.719) = 0.16670
Epoch: 08 -- f(0.725) = 0.15602
Epoch: 09 -- f(0.707) = 0.19025
Epoch: 10 -- f(0.692) = 0.21938
Epoch: 11 -- f(0.687) = 0.22943
Epoch: 12 -- f(0.689) = 0.22596
Epoch: 13 -- f(0.692) = 0.21937
Epoch: 14 -- f(0.694) = 0.21555
Epoch: 15 -- f(0.694) = 0.21499
Epoch: 16 -- f(0.694) = 0.21595
Epoch: 17 -- f(0.693) = 0.21690
Epoch: 18 -- f(0.693) = 0.21728
Epoch: 19 -- f(0.693) = 0.21723
```



Thanks!

Any questions?