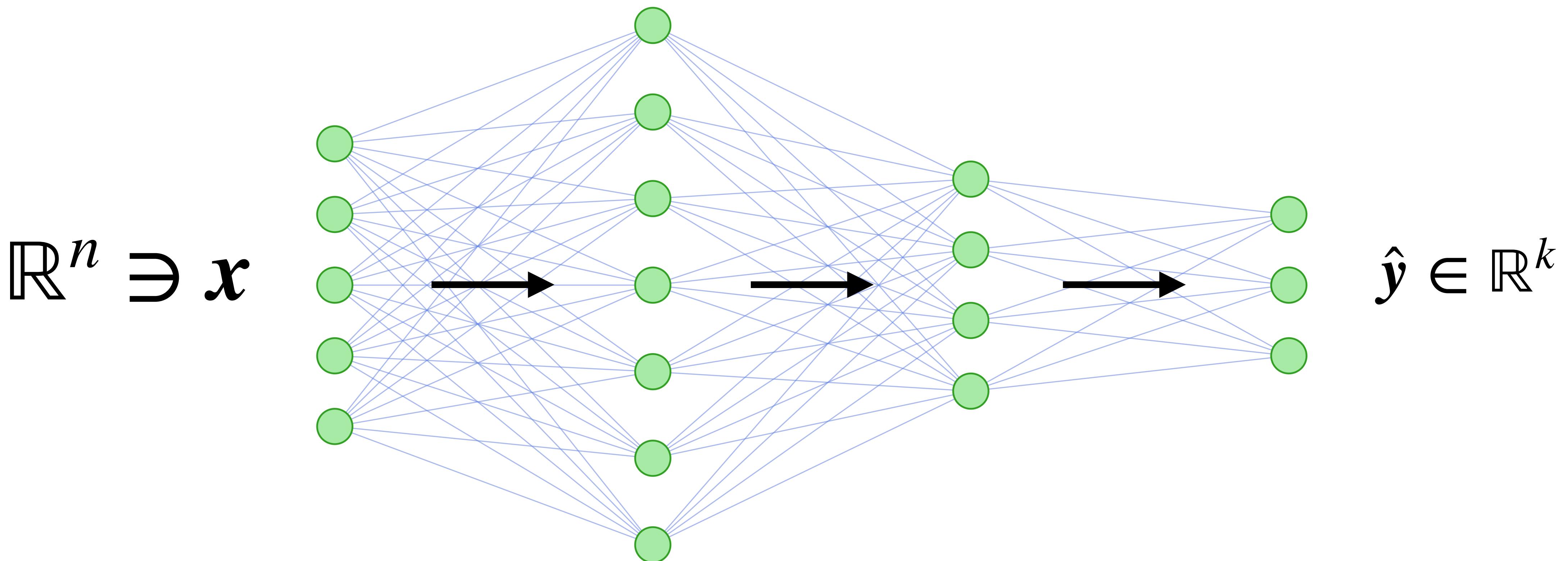


Neural Networks

Forward/backward propagation + training

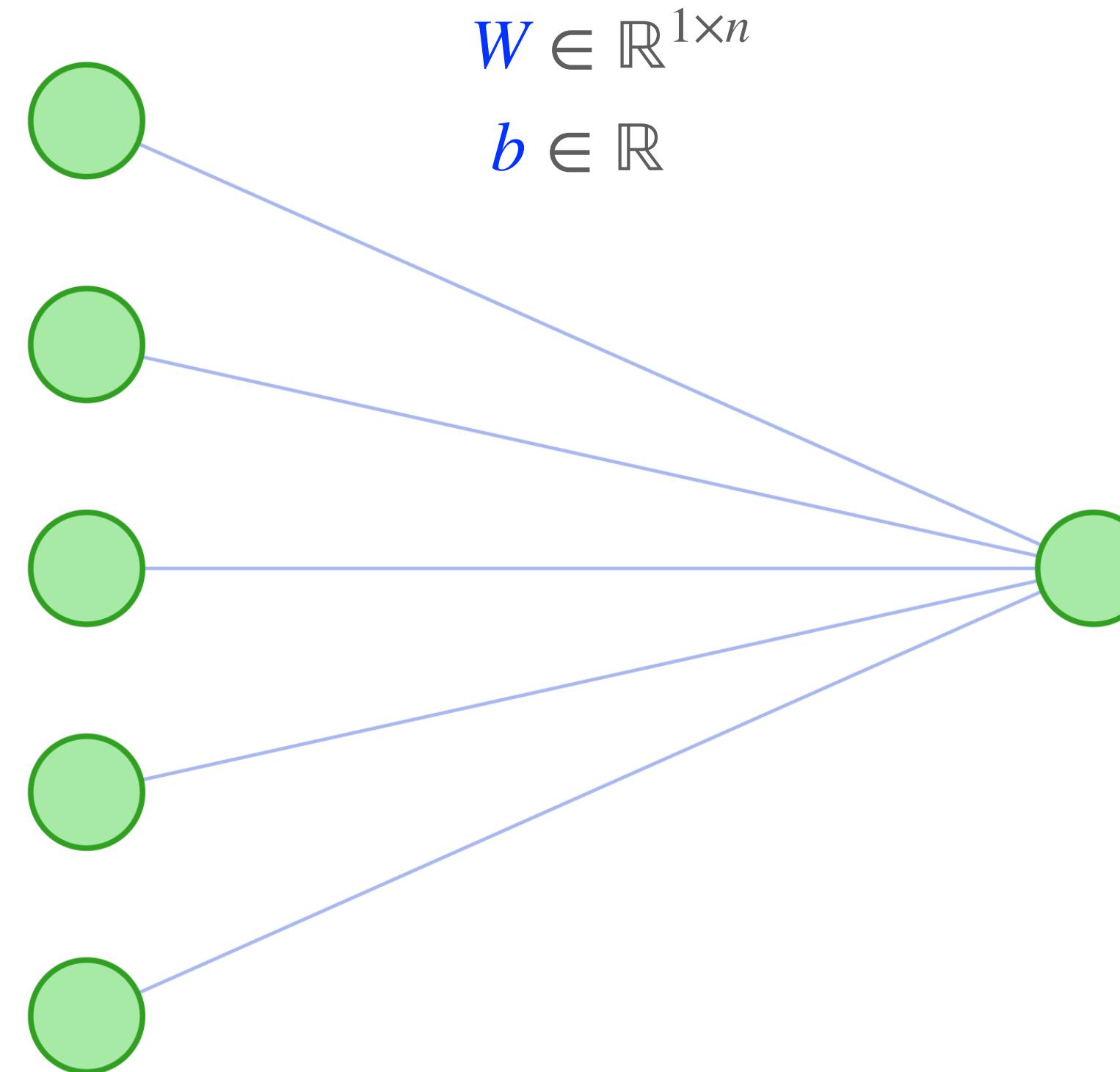
Forward propagation

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^k$$



Logistic regression (recap)

$\mathbb{R}^n \ni x$



$a \in (0, 1)$

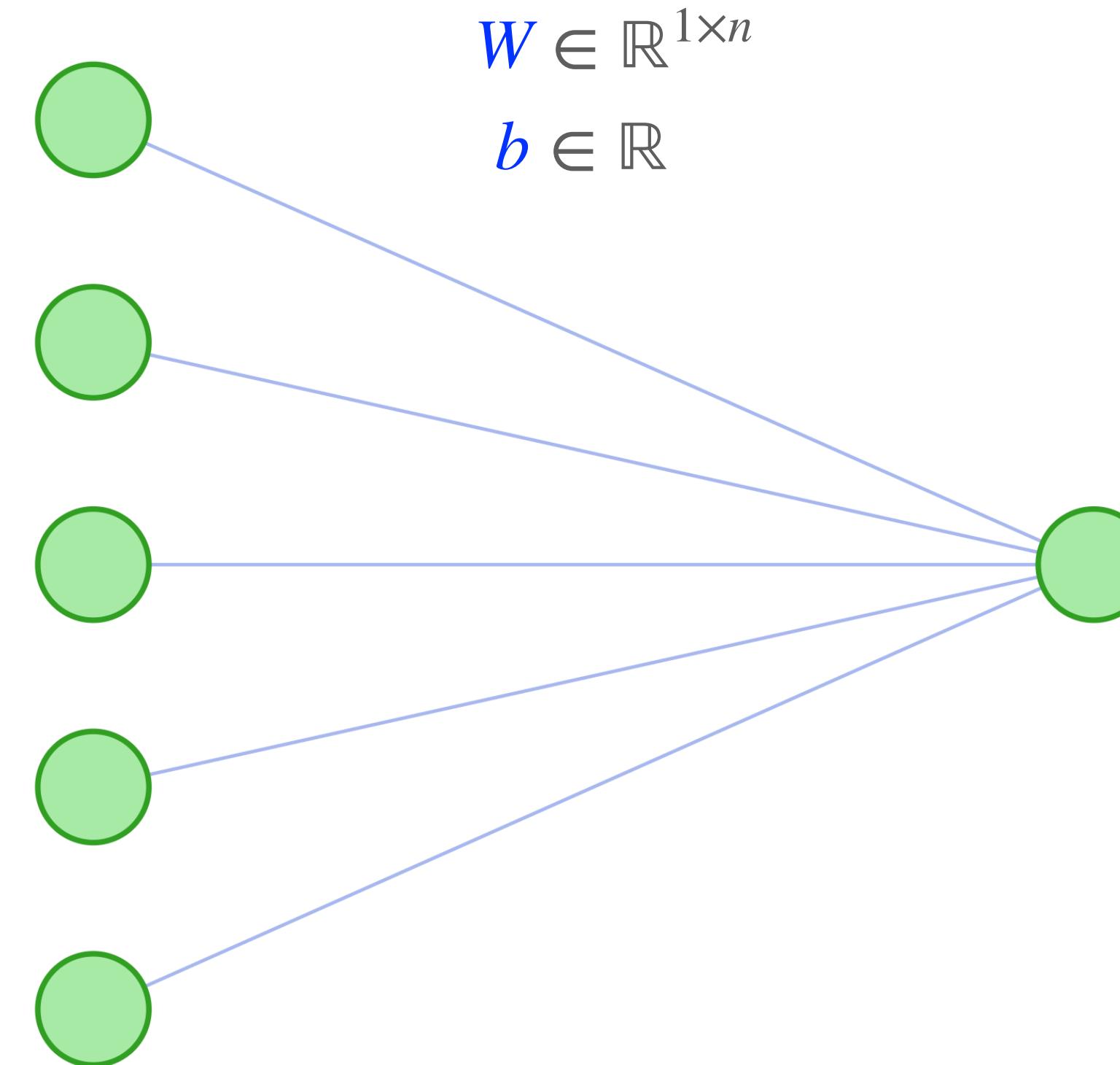
$$x^T \in \mathbb{R}^{1 \times n}$$

$$z = x^T W^T + b \in \mathbb{R}^{1 \times 1}$$

$$a = \sigma(z) \in \mathbb{R}^{1 \times 1}$$

Logistic regression (recap)

$\mathbb{R}^{N \times n} \ni X$



$a \in (0, 1)^N$

Efficient
implementation

$X \in \mathbb{R}^{N \times n}$

$$z = X W^T \oplus b \in \mathbb{R}^{N \times 1}$$

$$a = \sigma(z) \in \mathbb{R}^{N \times 1}$$

Forward propagation

Code demo in PyTorch



Forward propagation

```
In [1]: import torch  
  
torch.manual_seed(1); # for the reproducibility  
  
executed in 299ms, finished 17:46:14 2023-03-11
```

Linear model

```
In [2]: xT = torch.tensor([-8.0, 1.0, -3.0, 5.0, 2.0]).view(1, 5) # horizontal vector  
xT  
  
executed in 22ms, finished 17:46:14 2023-03-11  
  
Out[2]: tensor([-8., 1., -3., 5., 2.])
```

```
In [3]: xT.shape  
  
executed in 19ms, finished 17:46:14 2023-03-11  
  
Out[3]: torch.Size([1, 5])
```

```
In [4]: W = torch.rand(5).view(1, 5) # random numbers from (0, 1)  
b = torch.rand(1)  
  
executed in 2ms, finished 17:46:14 2023-03-11
```

```
In [5]: W  
  
executed in 2ms, finished 17:46:14 2023-03-11  
  
Out[5]: tensor([0.7576, 0.2793, 0.4031, 0.7347, 0.0293])
```

```
In [6]: b  
  
executed in 3ms, finished 17:46:14 2023-03-11  
  
Out[6]: tensor([0.7999])
```

```
In [7]: W.shape, b.shape  
  
executed in 2ms, finished 17:46:14 2023-03-11  
  
Out[7]: (torch.Size([1, 5]), torch.Size([1]))
```

$$z = \mathbf{x}^T \mathbf{W}^T + \mathbf{b} \in \mathbb{R}^{1 \times 1}$$

```
In [8]: z = torch.mm(xT, W.T) + b  
z.shape  
  
executed in 2ms, finished 17:46:14 2023-03-11  
  
Out[8]: torch.Size([1, 1])
```

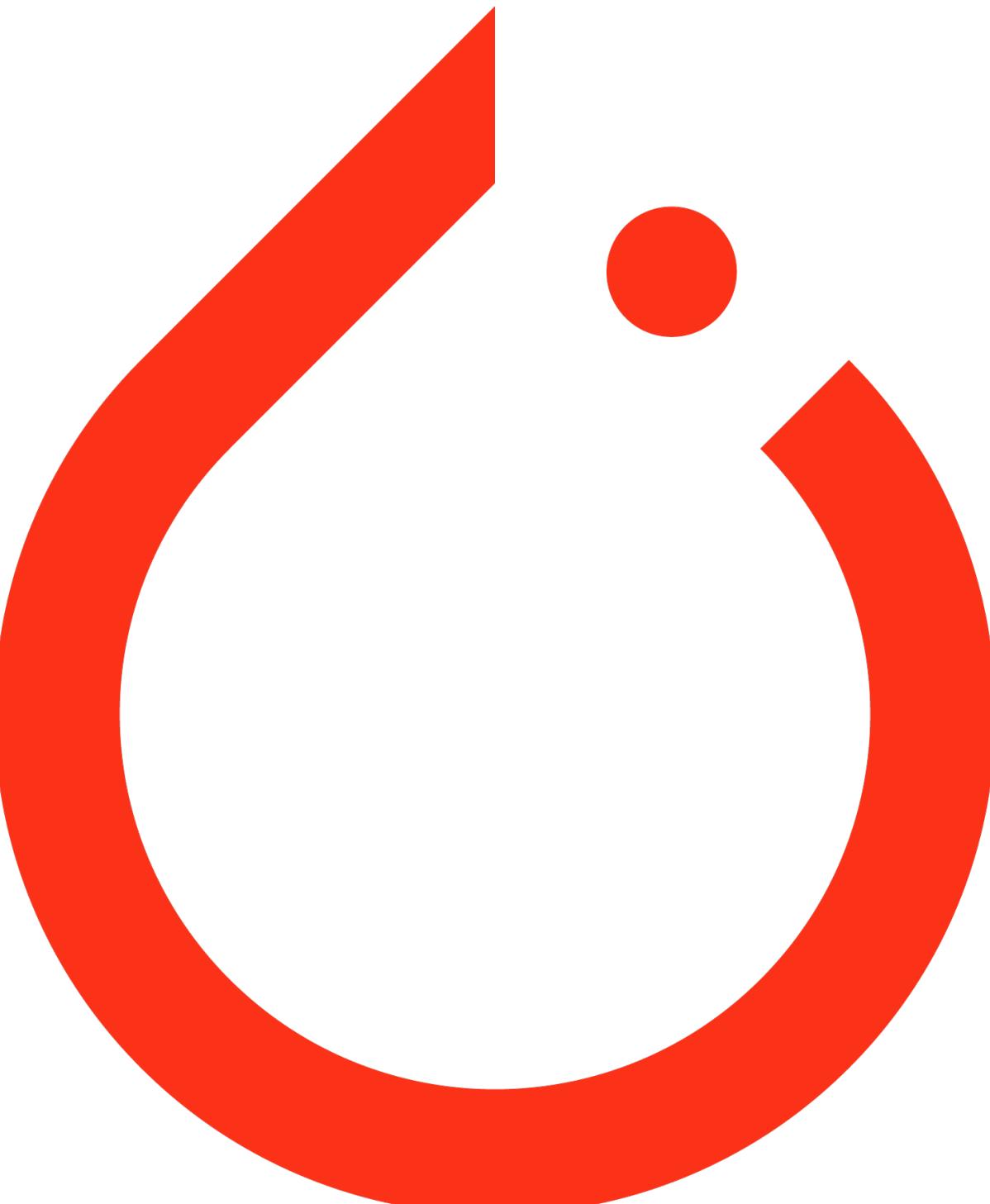
```
In [9]: z  
  
executed in 3ms, finished 17:46:14 2023-03-11  
  
Out[9]: tensor([-2.4591])
```

$$a = \sigma(z) \in \mathbb{R}^{1 \times 1}$$

```
In [10]: a = torch.sigmoid(z)  
a  
  
executed in 3ms, finished 17:46:14 2023-03-11  
  
Out[10]: tensor([0.0788])
```

Forward propagation

Code demo in PyTorch



Forward propagation

```
In [1]: import torch  
torch.manual_seed(1); # for the reproducibility  
executed in 284ms, finished 17:43:29 2023-03-11
```

Linear model

```
In [2]: X = torch.tensor([[-8.0, 1.0, -3.0, 5.0, 2.0], [-2.0, 2.0, -1.0, 5.0, 2.0], [7.0, 1.0, 9.0, -5.0, -2.0]])  
X  
executed in 17ms, finished 17:43:29 2023-03-11  
Out[2]: tensor([-8., 1., -3., 5., 2.],  
[-2., 2., -1., 5., 2.],  
[ 7., 1., 9., -5., -2.])
```

```
In [3]: X.shape  
executed in 25ms, finished 17:43:29 2023-03-11
```

```
Out[3]: torch.Size([3, 5])
```

```
In [4]: W = torch.rand(5).view(1, 5) # random numbers from (0, 1)  
b = torch.rand(1)  
executed in 8ms, finished 17:43:29 2023-03-11
```

```
In [5]: W  
executed in 2ms, finished 17:43:29 2023-03-11
```

```
Out[5]: tensor([0.7576, 0.2793, 0.4031, 0.7347, 0.0293])
```

```
In [6]: b  
executed in 3ms, finished 17:43:29 2023-03-11
```

```
Out[6]: tensor(0.7999)
```

```
In [7]: W.shape, b.shape  
executed in 2ms, finished 17:43:29 2023-03-11
```

```
Out[7]: (torch.Size([1, 5]), torch.Size([1]))
```

$$z = XW^T \oplus b \in \mathbb{R}^{N \times 1}$$

```
In [8]: z = torch.mm(X, W.T) + b  
z.shape  
executed in 2ms, finished 17:43:29 2023-03-11
```

```
Out[8]: torch.Size([3, 1])
```

```
In [9]: z  
executed in 2ms, finished 17:43:29 2023-03-11
```

```
Out[9]: tensor([-2.4591,  
[ 3.1721],  
[ 6.2782]])
```

$$a = \sigma(z) \in \mathbb{R}^{N \times 1}$$

```
In [10]: a = torch.sigmoid(z)  
a  
executed in 2ms, finished 17:43:29 2023-03-11
```

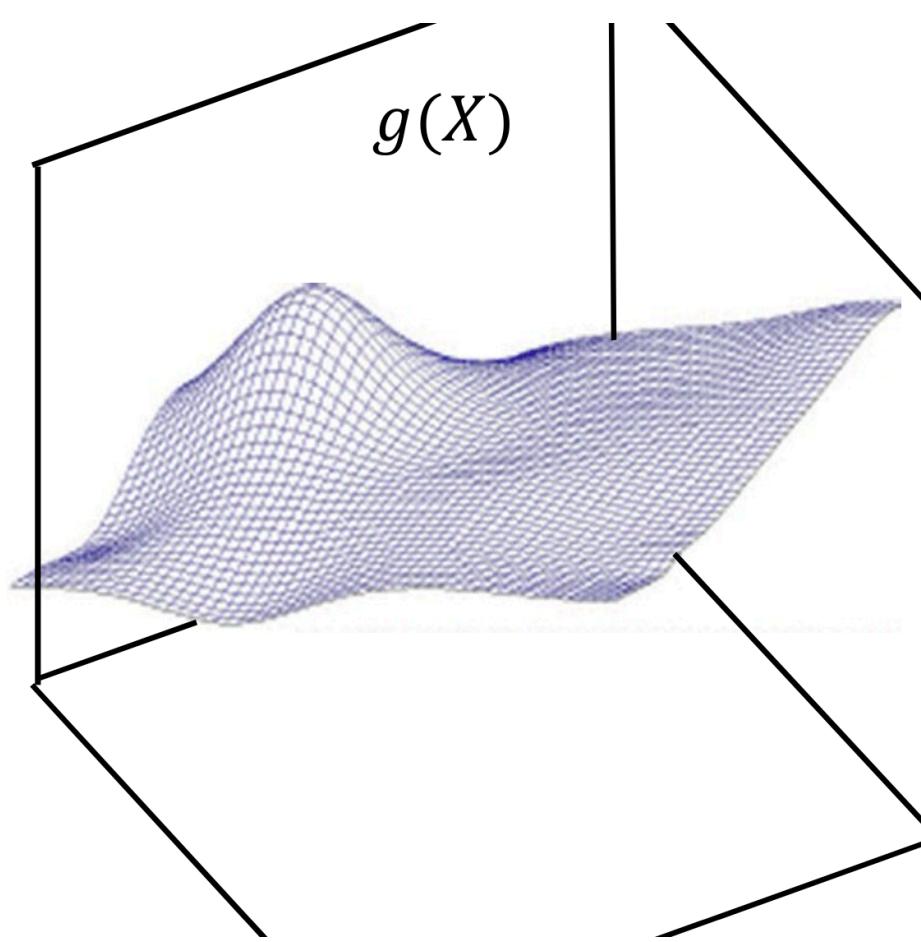
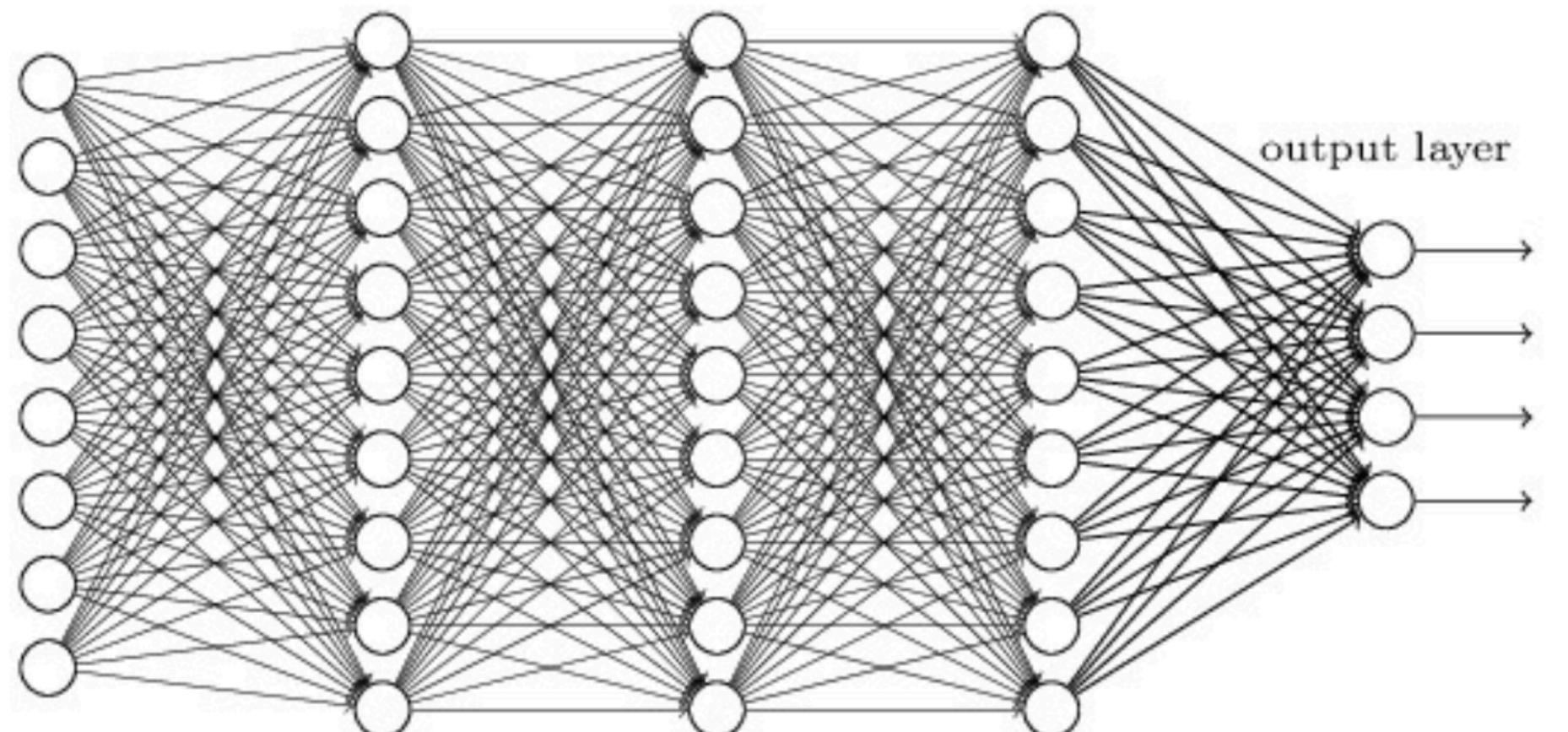
```
Out[10]: tensor([[0.0788],  
[0.9598],  
[0.9981]])
```

Supervised learning

- regression
- classification
- ...

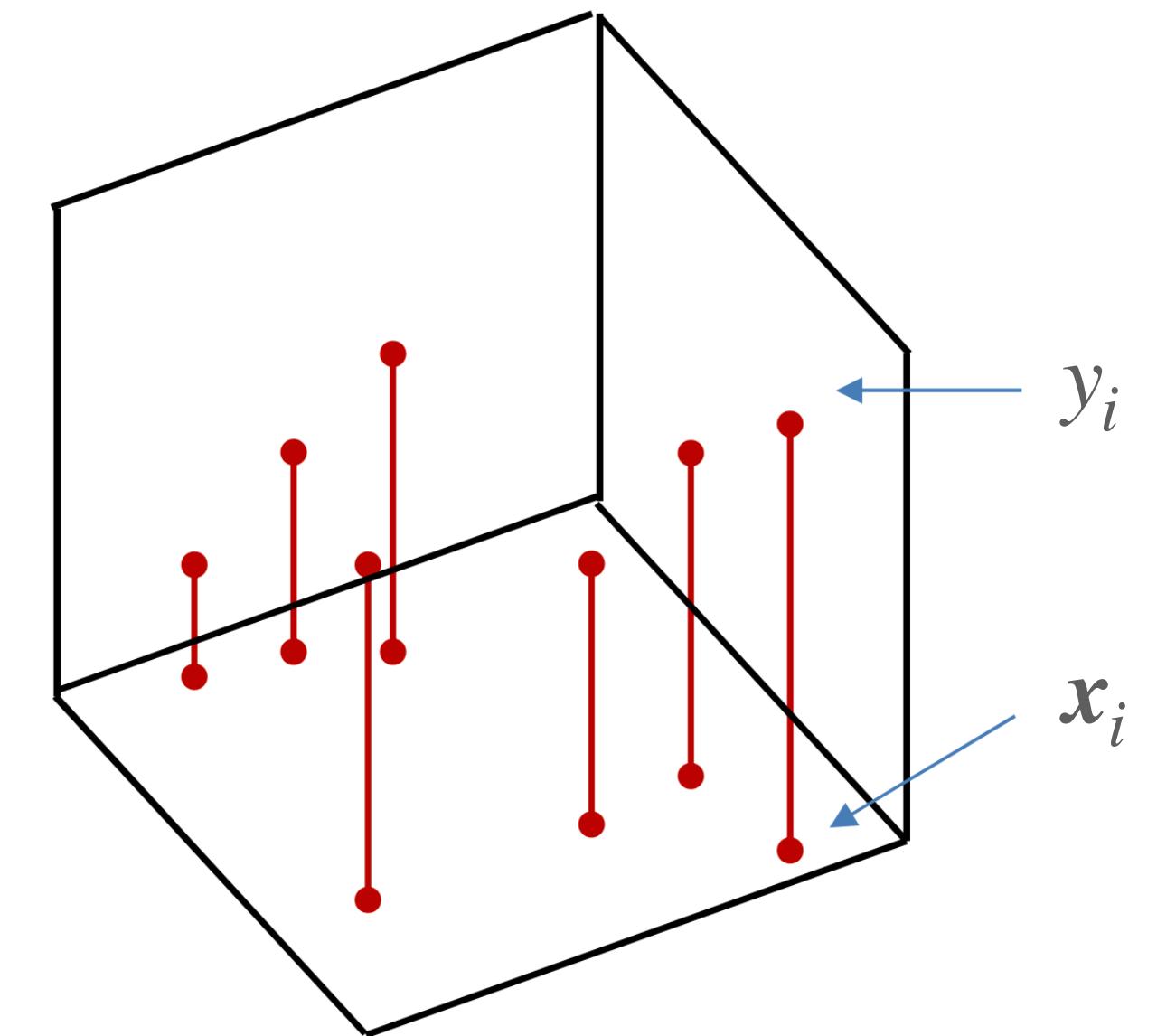
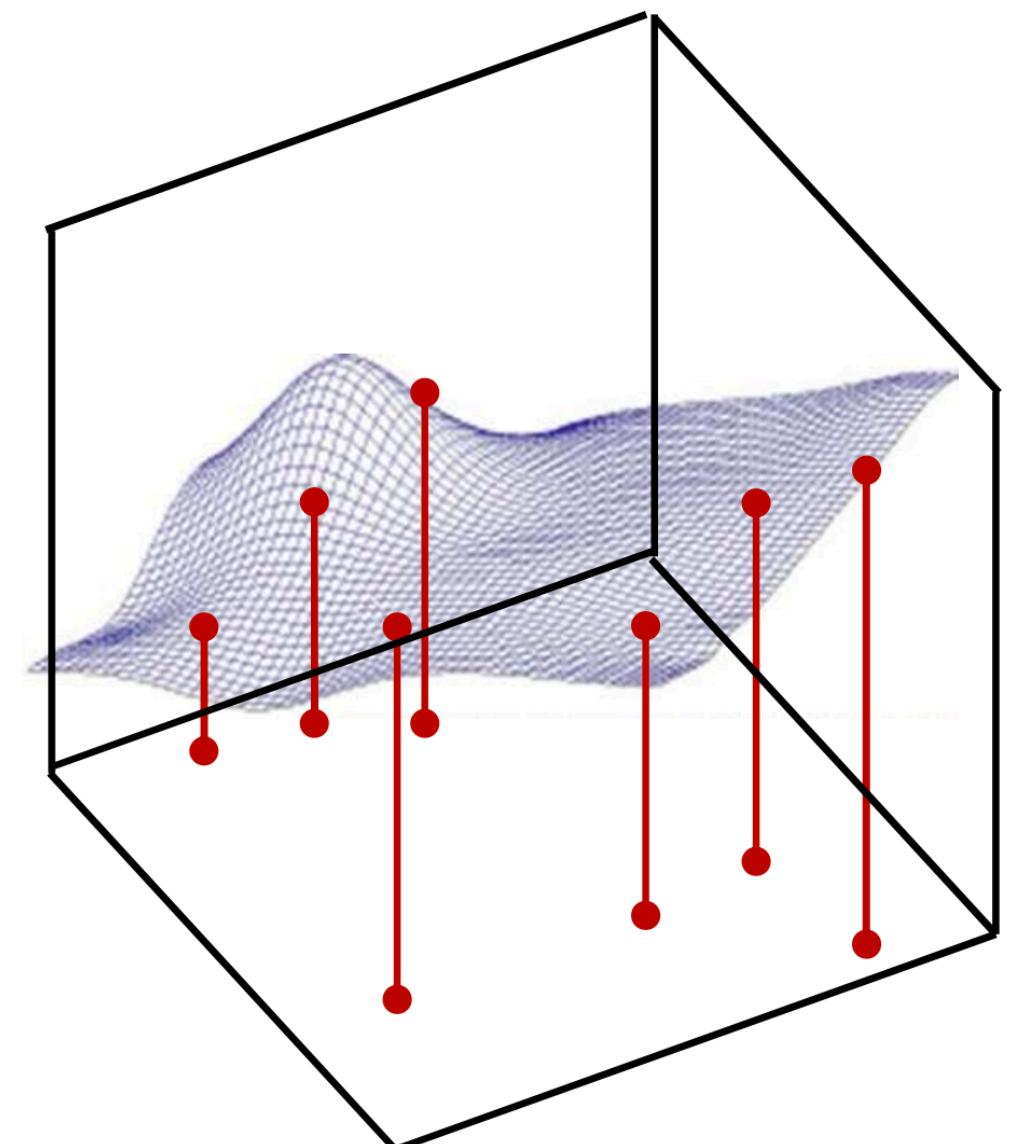
What we learn?

- network parameters W
- $f(X, W) \approx g(X)$



Learning by sampling

- training set $X_{\text{train}} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T \in \mathbb{R}^{N \times n}$, $y_{\text{train}} = [y_1, y_2, \dots, y_N]^T \in \mathbb{R}^{N \times k}$
- $f(X, \mathbf{W}) \approx g(X)$
- Questions:
 - good sampling?
- Examples:
 - images + labels
 - speech recordings and their transcription



Training the neural network

- Find the parameters \mathbf{W} that minimize the total average **loss**
- Total average loss:

$$\text{LOSS} = \frac{1}{N} \sum_{p=1}^N \textcolor{red}{div}(\hat{\mathbf{y}}_p, \mathbf{y}_p) \quad \hat{\mathbf{y}}_p := f(\mathbf{x}_p ; \mathbf{W})$$

Output Activation + Loss Function

- Output activation:
 - regression -> **identity** (id)
 - classification
 - binary -> **sigmoid**
 - multiclass -> **softmax**
- Loss function
 - regression → mean squared error (MSE)
 - classification → cross entropy (CE)

Linear regression

MSE Loss function

- (scaled) L_2 error (squared euclidean distance)

$$f(\mathbf{x}; \mathbf{w}, b) = \hat{y}$$

$$\hat{y} = \mathbf{x}\mathbf{w}^T + b$$

$$\text{div}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$$

Logistic regression

Cross entropy

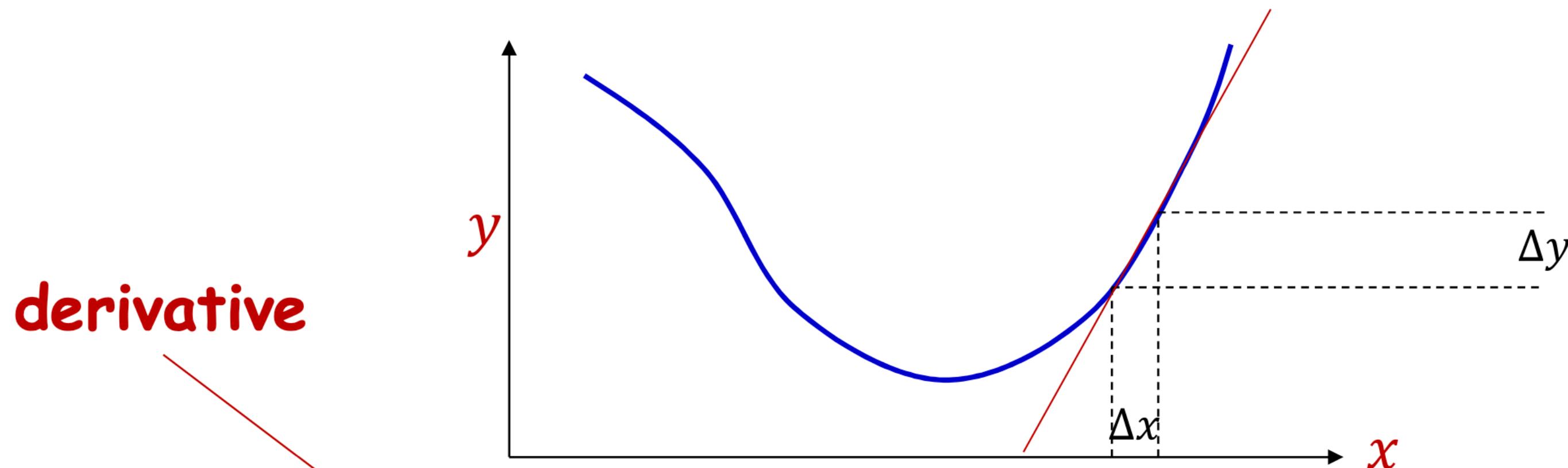
- binary cross entropy loss

$$f(\mathbf{x}; \mathbf{w}, b) = \hat{y}$$

$$\hat{y} = \sigma(\mathbf{x}\mathbf{w}^T + b)$$

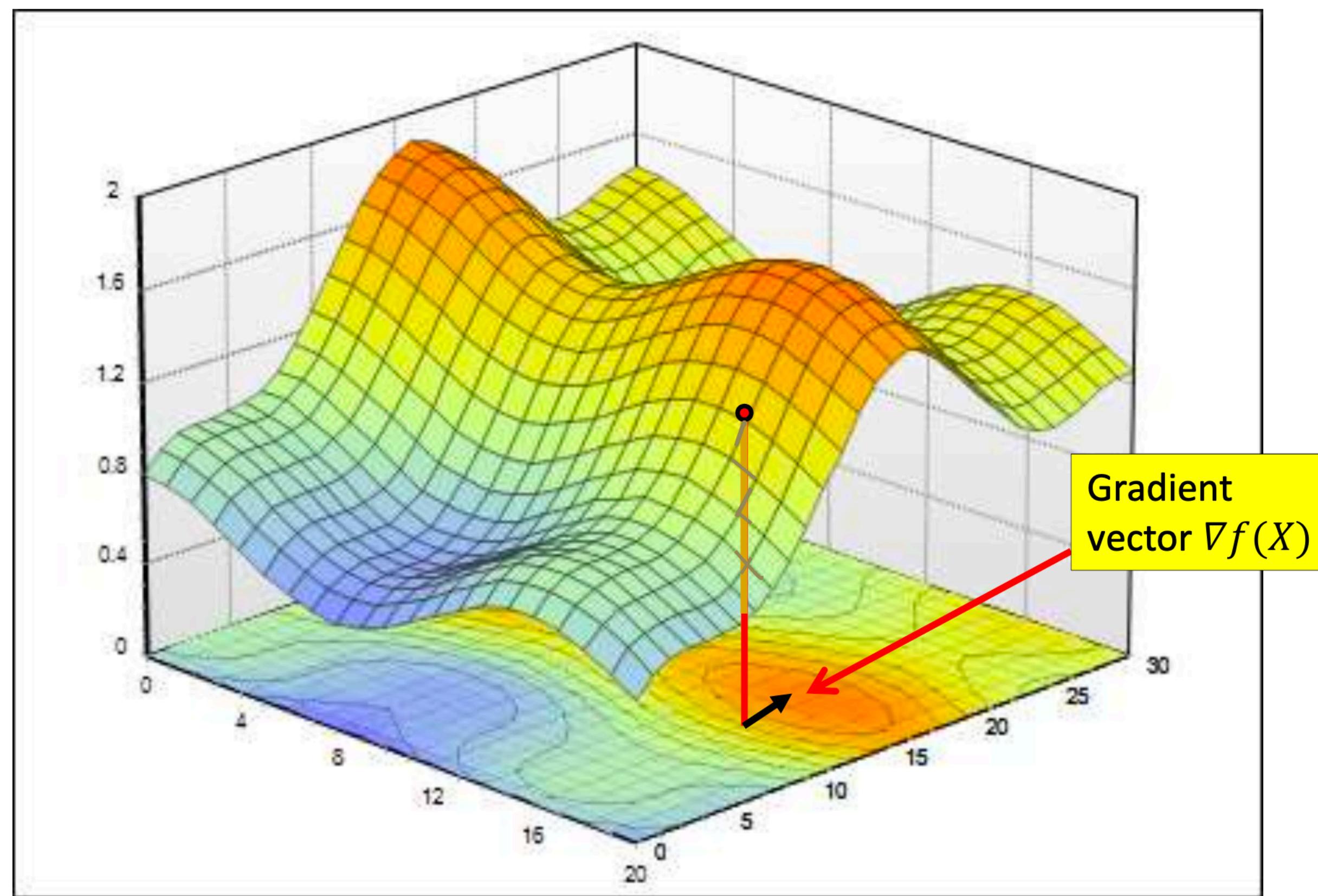
$$\text{div}(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

A brief note on derivatives..

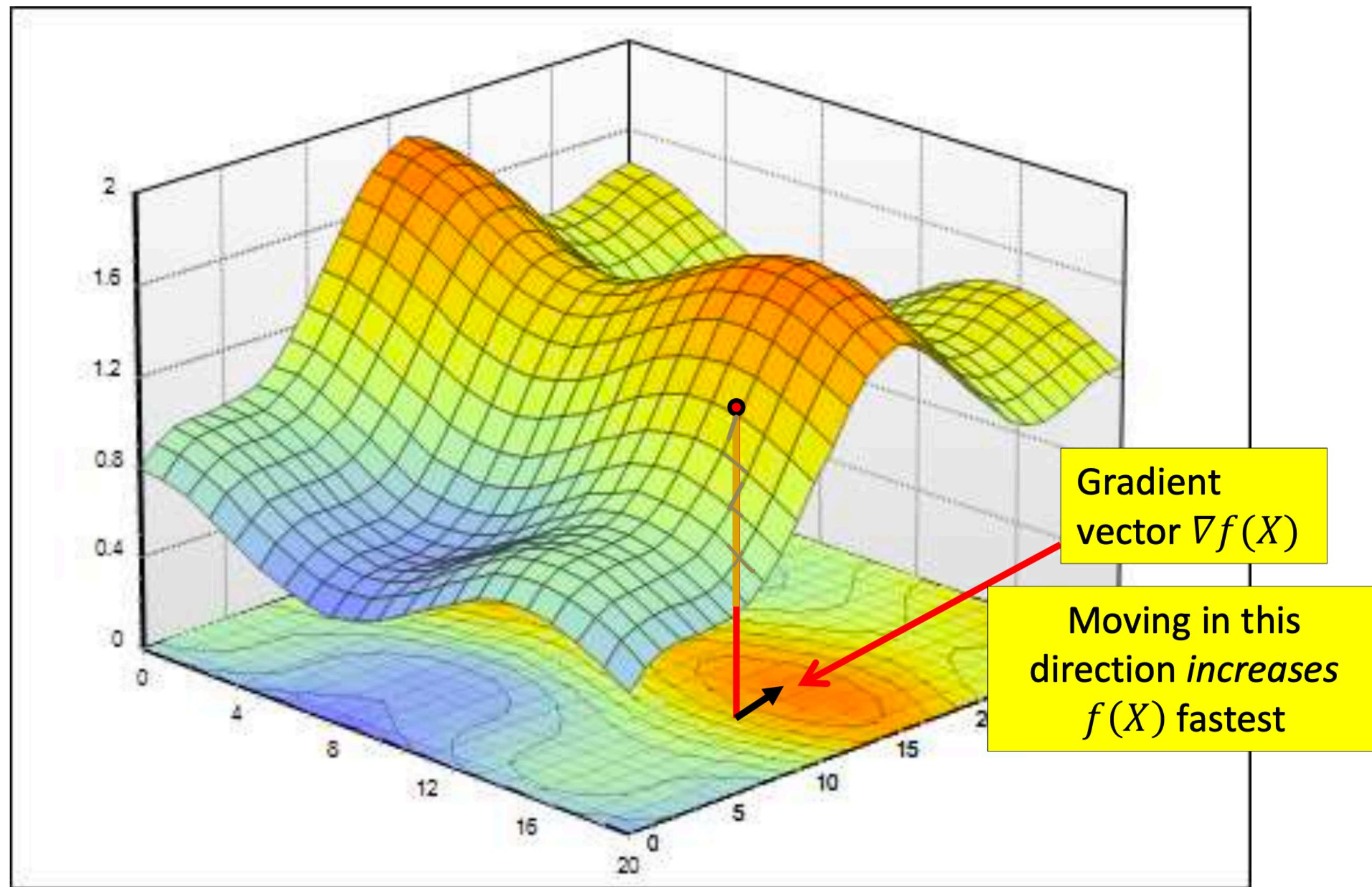


- A derivative of a function at any point tells us how much a minute increment to the *argument* of the function will increment the *value* of the function
 - For any $y = f(x)$, expressed as a multiplier α to a tiny increment Δx to obtain the increments Δy to the output
$$\Delta y = \alpha \Delta x$$
 - Based on the fact that at a fine enough resolution, any smooth, continuous function is locally linear at any point

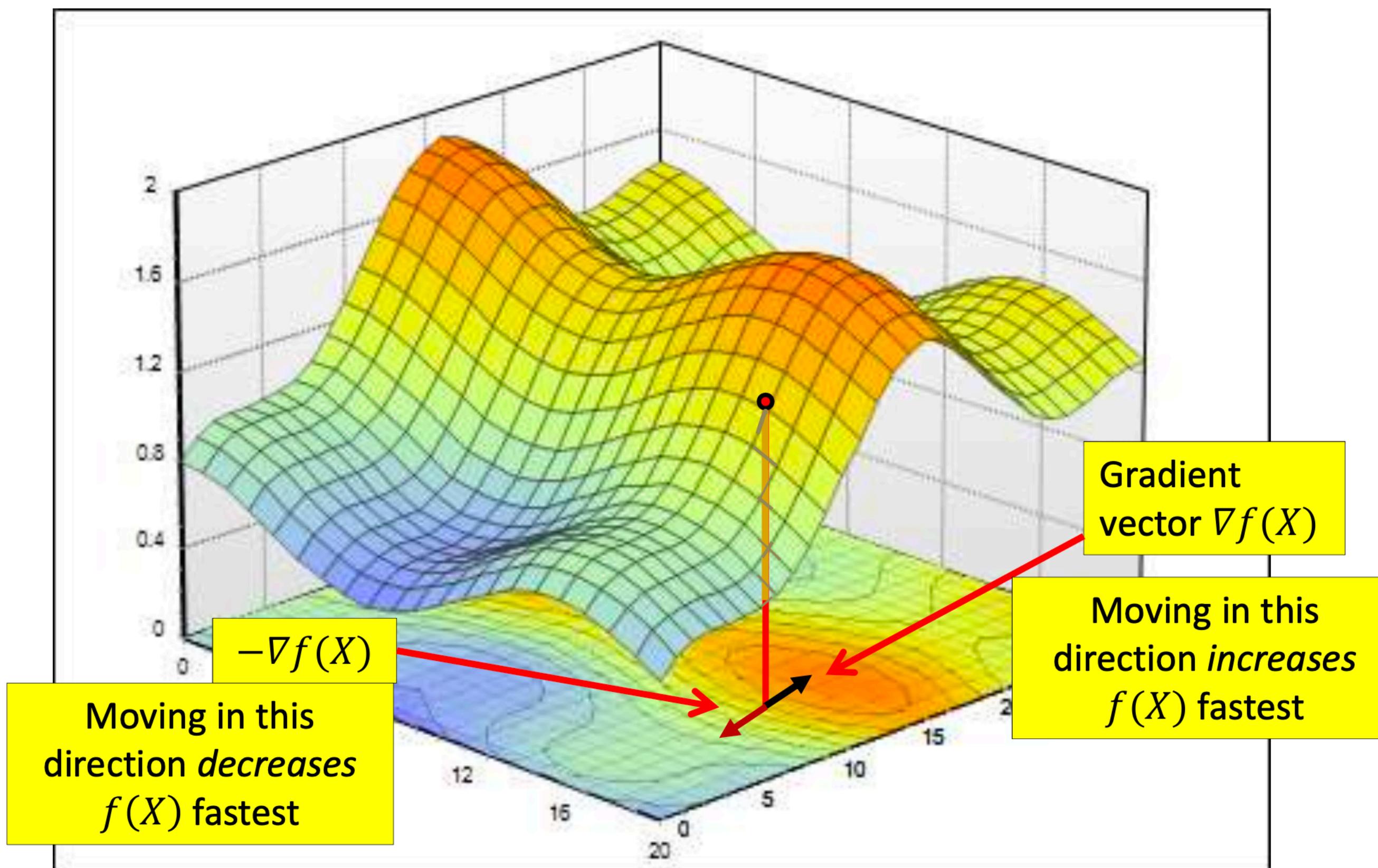
Gradient



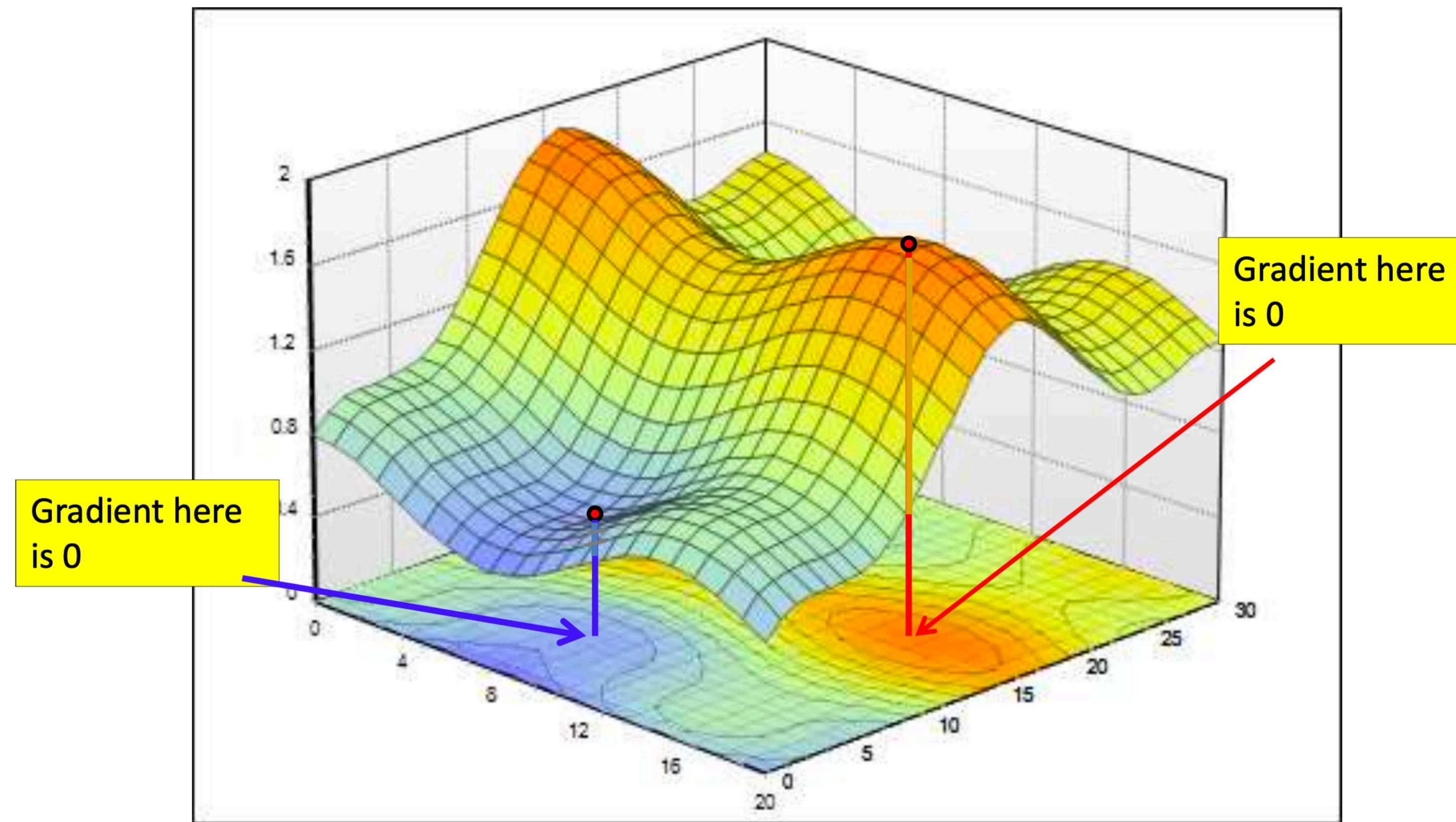
Gradient

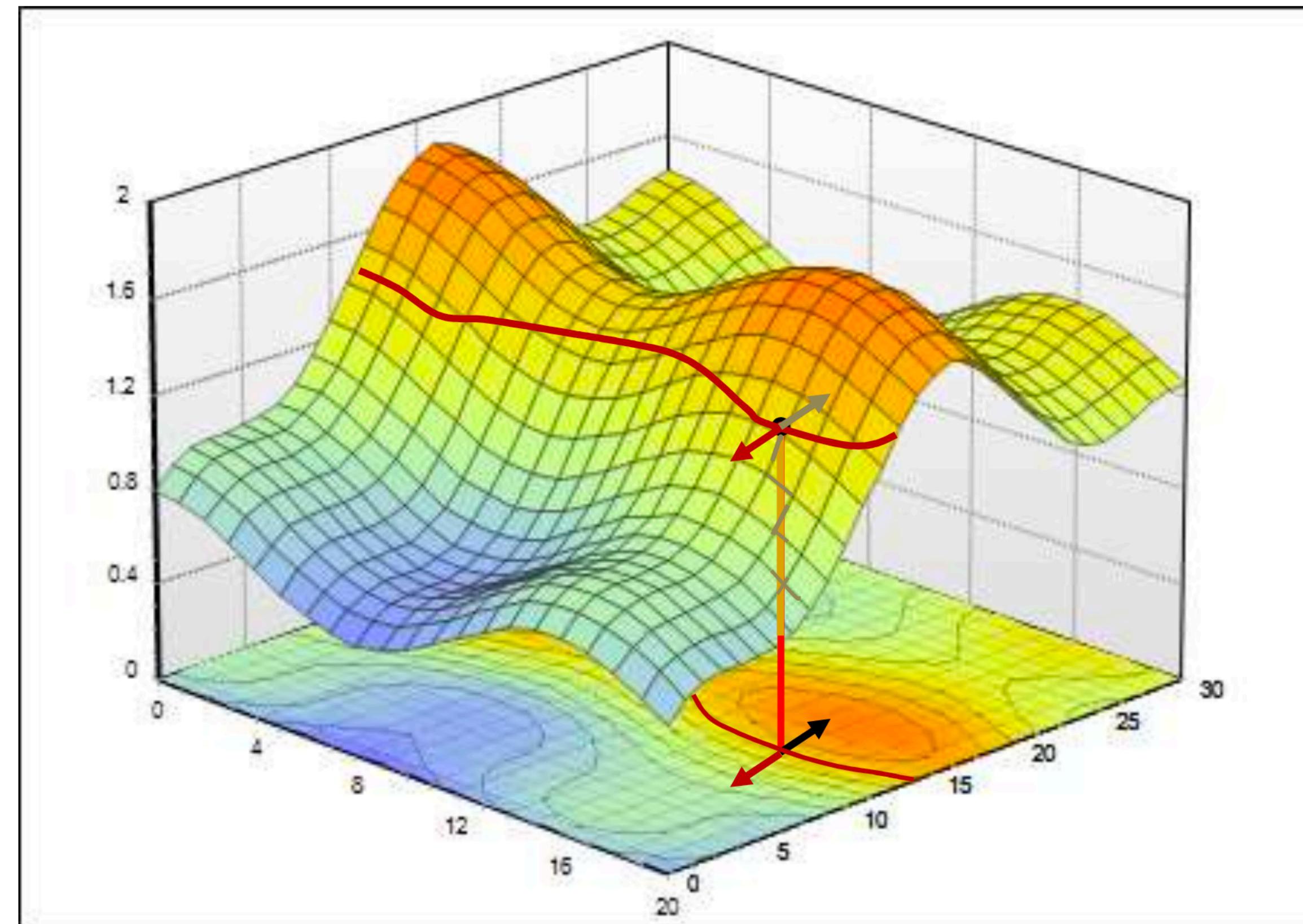


Gradient



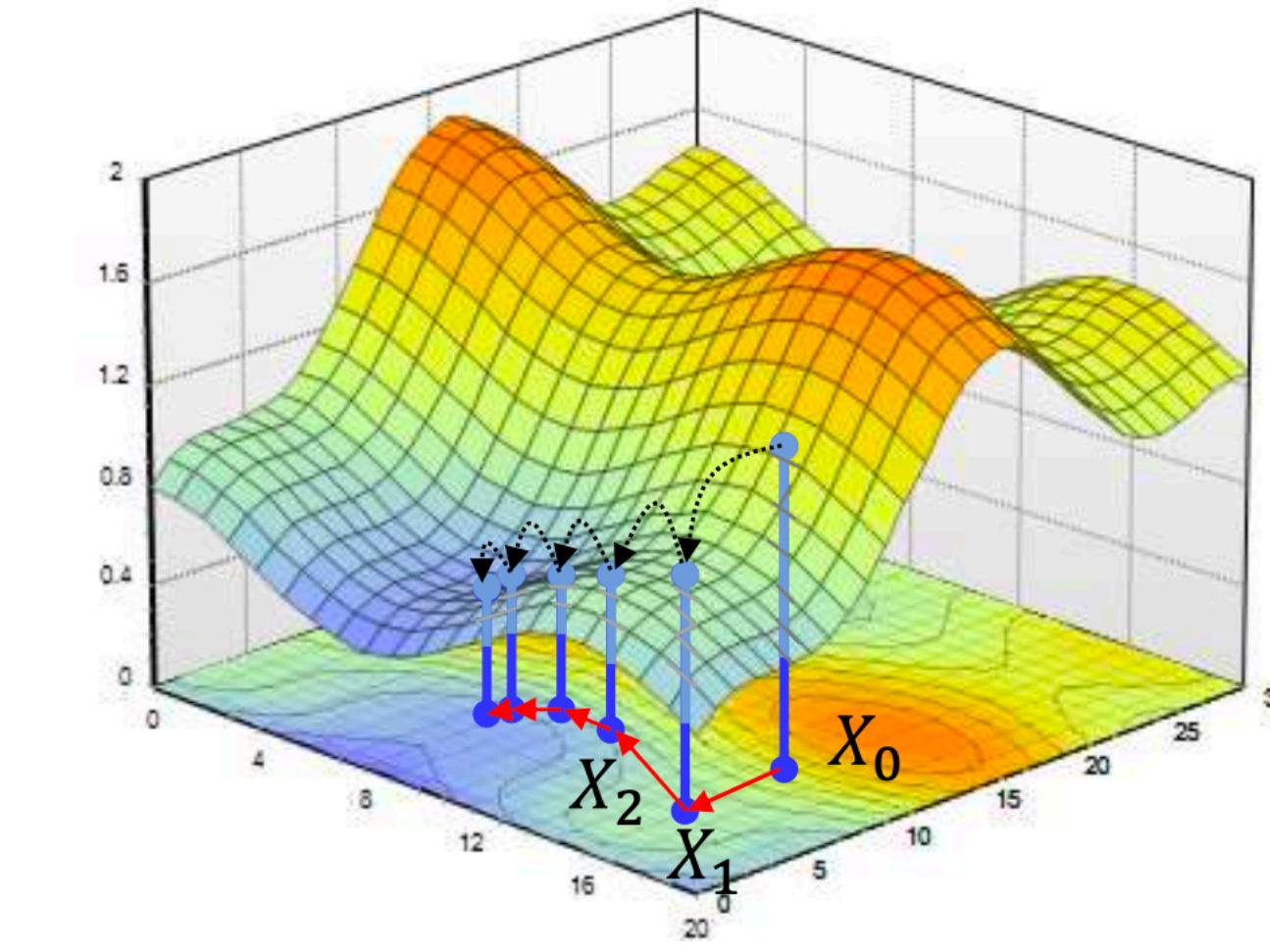
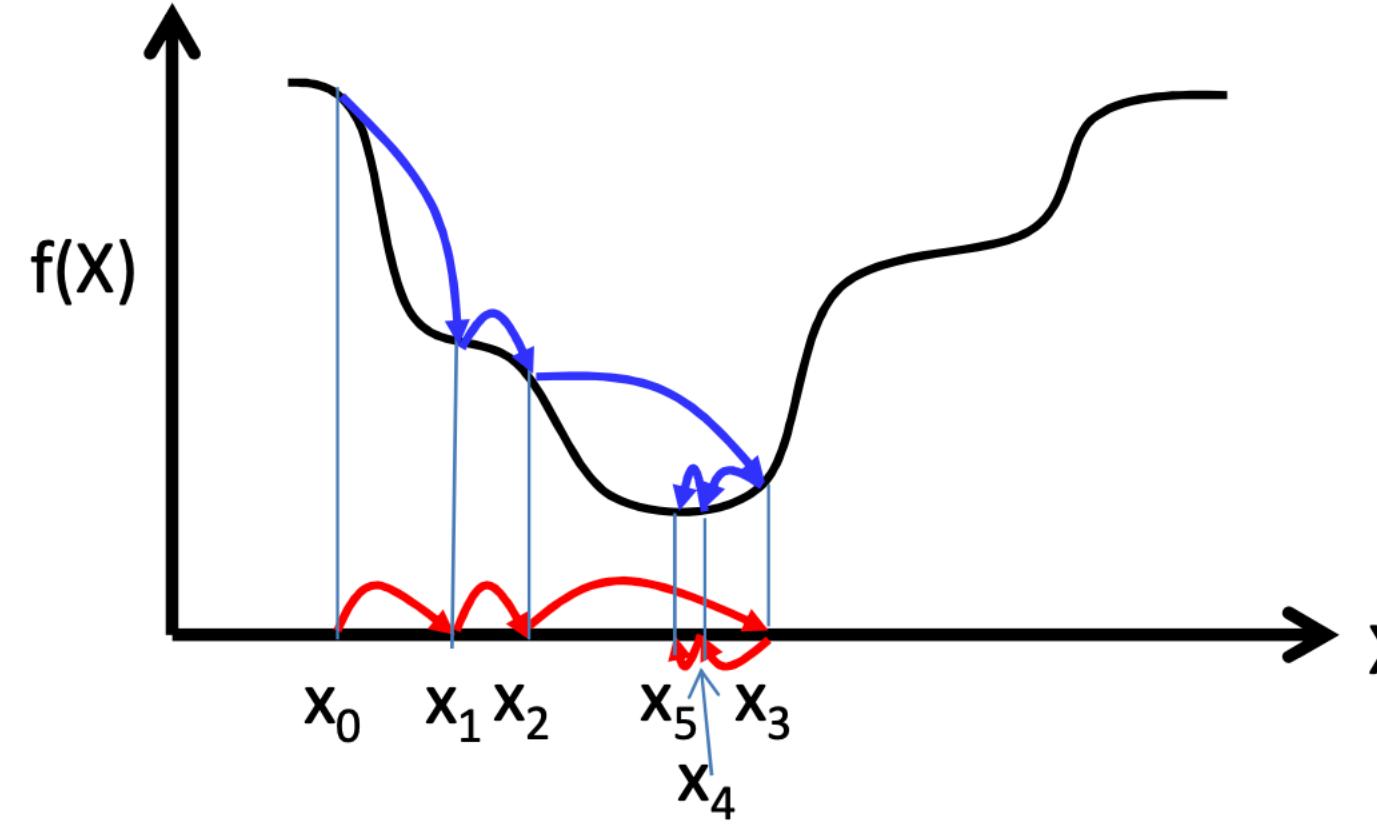
Gradient





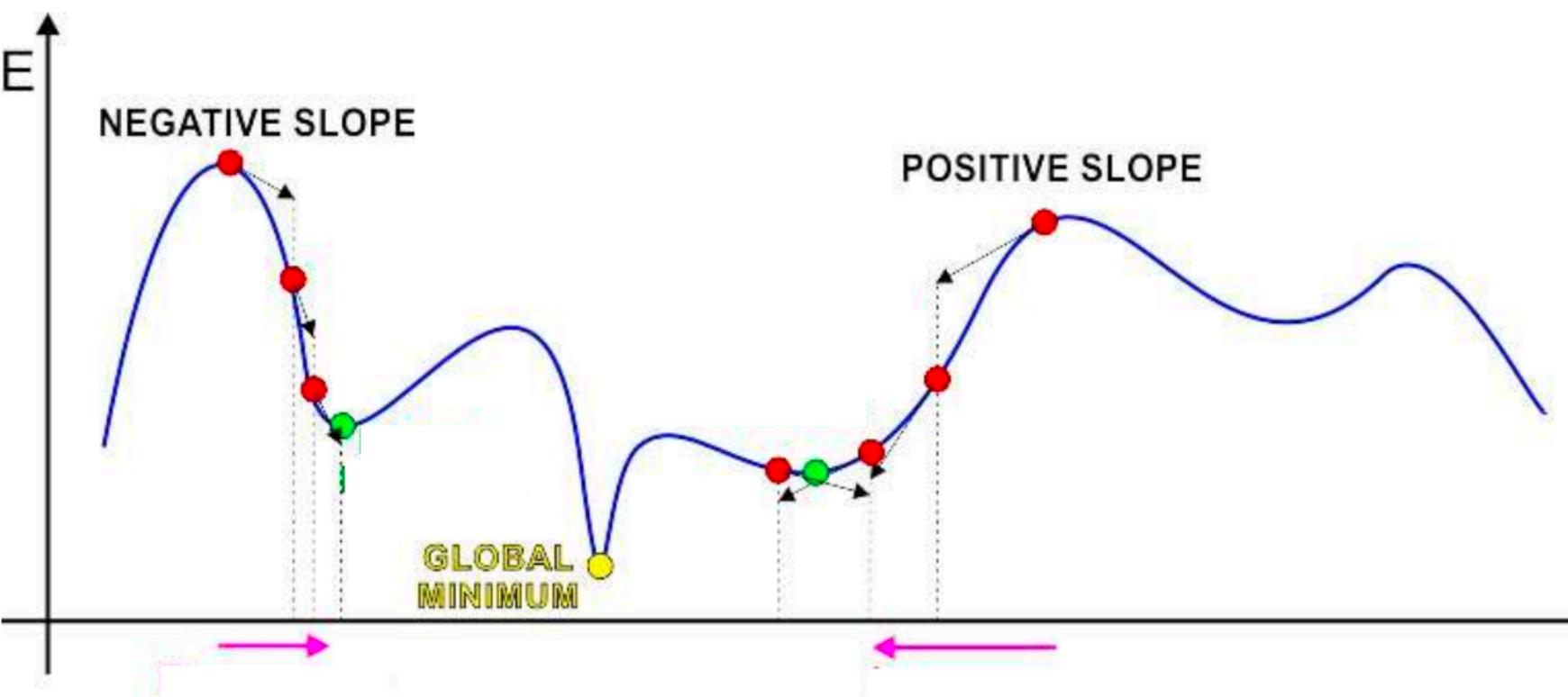
- The gradient vector $\nabla f(X)$ is perpendicular to the level curve

Iterative solutions



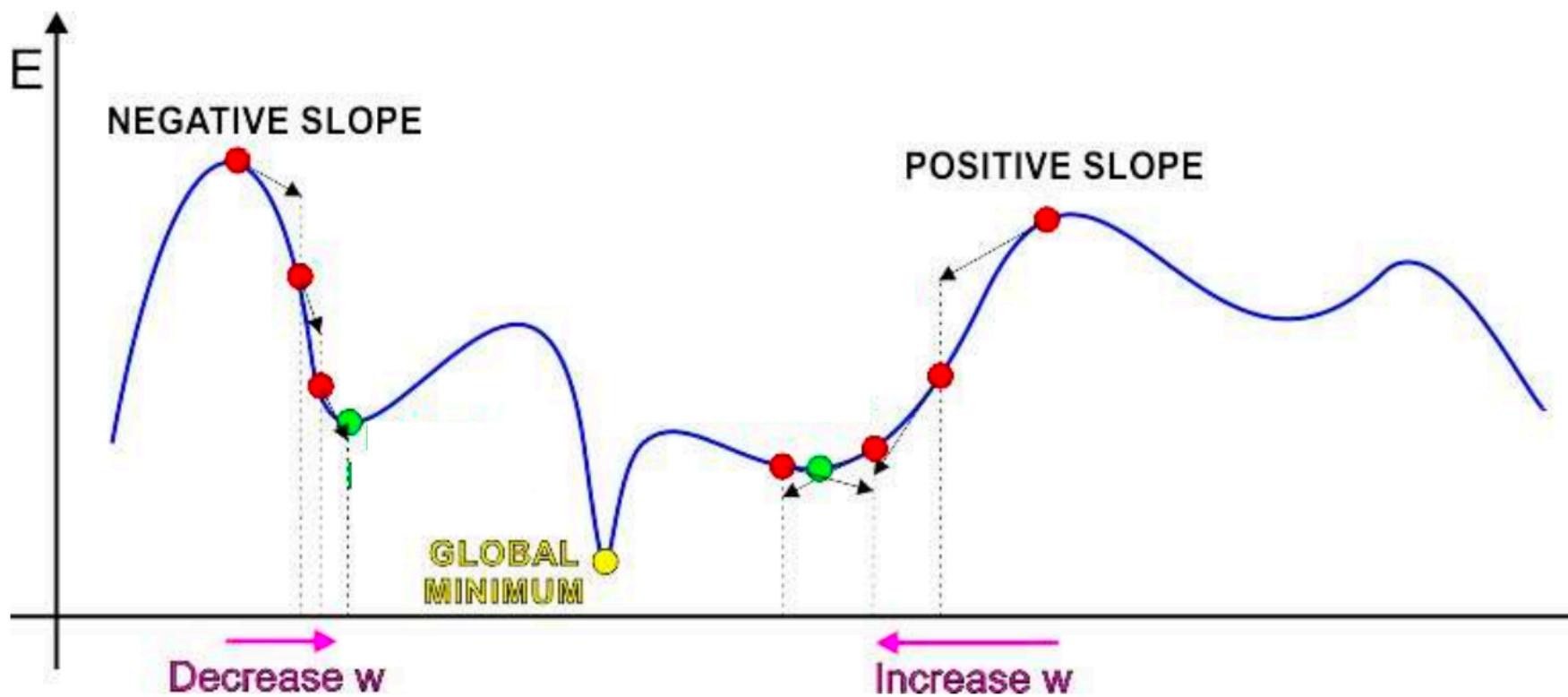
- Iterative solutions
 - Start from an initial guess X_0 for the optimal X
 - Update the guess towards a (hopefully) “better” value of $f(X)$
 - Stop when $f(X)$ no longer decreases
- Problems:
 - Which direction to step in
 - How big must the steps be

The Approach of Gradient Descent



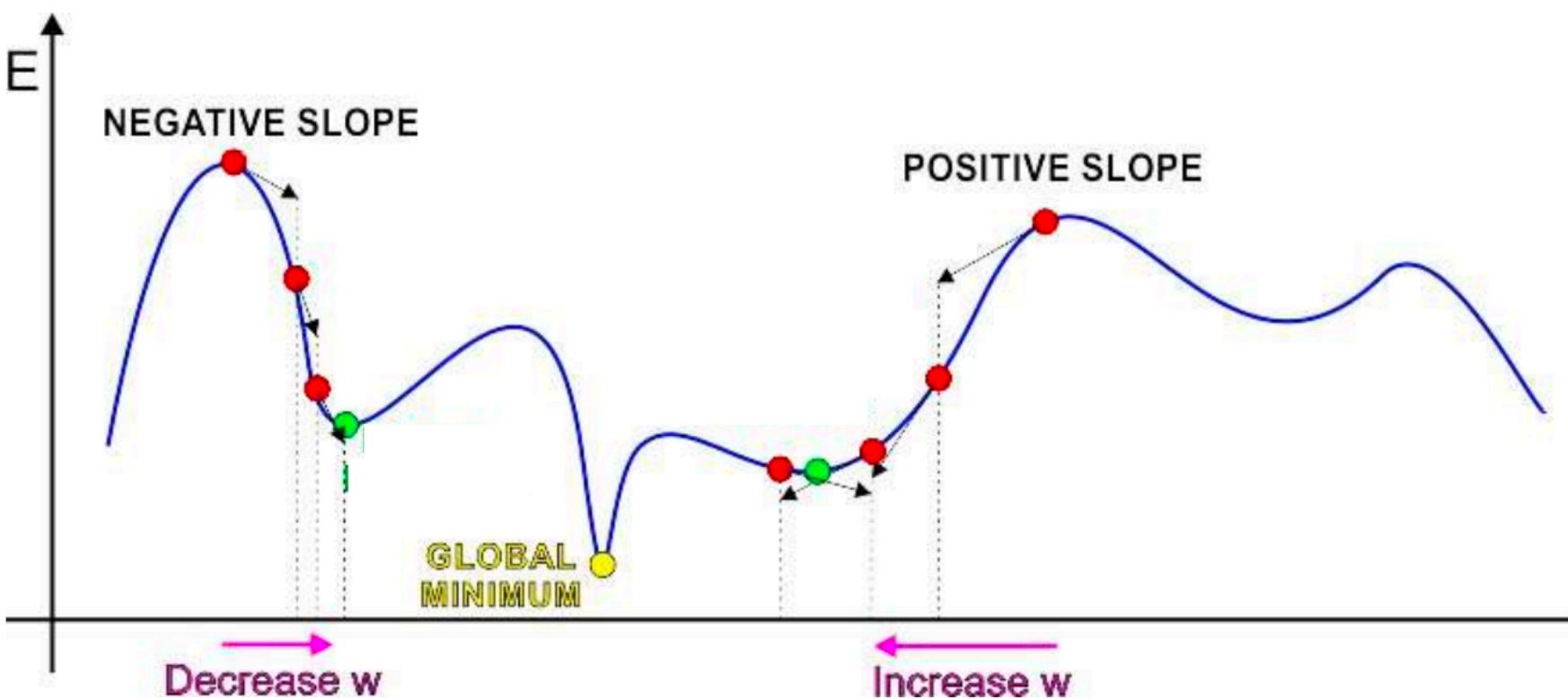
- Iterative solution:
 - Start at some point
 - Find direction in which to shift this point to decrease error
 - This can be found from the derivative of the function
 - A positive derivative → moving left decreases error
 - A negative derivative → moving right decreases error
 - Shift point in this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

The Approach of Gradient Descent



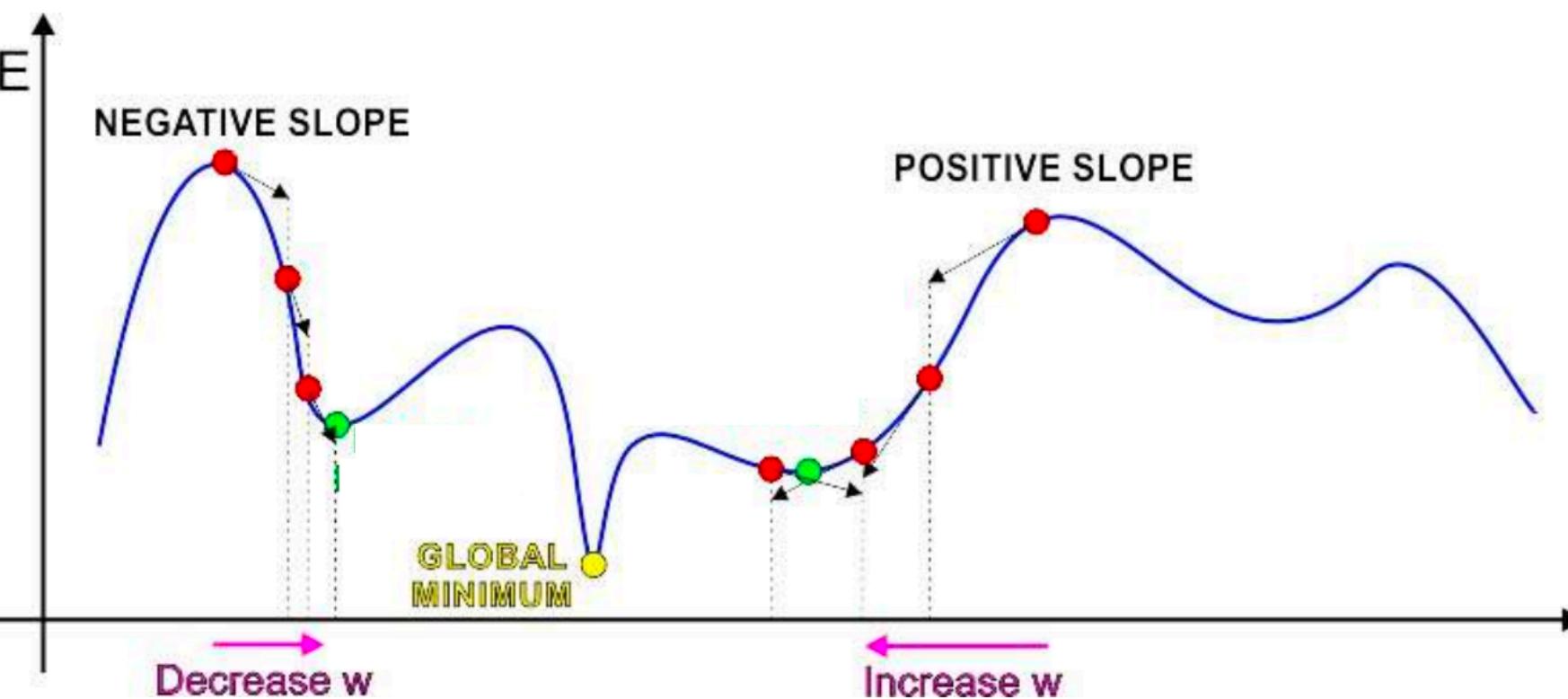
- Iterative solution: Trivial algorithm
 - Initialize x^0
 - While $f'(x^k) \neq 0$
 - If $\text{sign}(f'(x^k))$ is positive:
$$x^{k+1} = x^k - \text{step}$$
 - Else
$$x^{k+1} = x^k + \text{step}$$
 - What must step be to ensure we actually get to the optimum?

The Approach of Gradient Descent



- Iterative solution: Trivial algorithm
 - Initialize x^0
 - While $f'(x^k) \neq 0$
$$x^{k+1} = x^k - sign(f'(x^k)).step$$
- Identical to previous algorithm

The Approach of Gradient Descent



- Iterative solution: Trivial algorithm
 - Initialize x^0
 - While $f'(x^k) \neq 0$
$$x^{k+1} = x^k - \eta^k f'(x^k)$$
- η^k is the “step size”

Gradient descent/ascent (multivariate)

- The gradient descent/ascent method to find the minimum or maximum of a function f iteratively

- To find a *maximum* move *in the direction of the gradient*

$$x^{k+1} = x^k + \eta^k \nabla f(x^k)^T$$

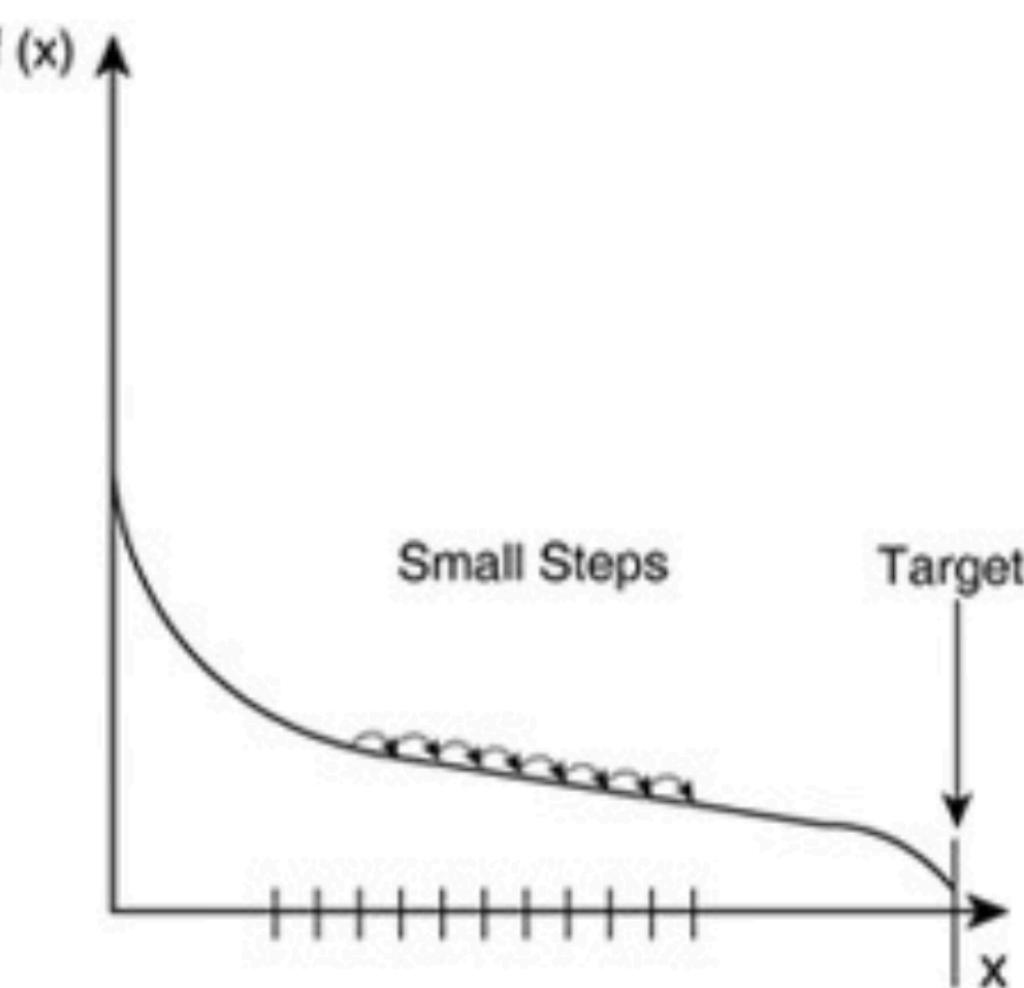
- To find a *minimum* move *exactly opposite the direction of the gradient*

$$x^{k+1} = x^k - \eta^k \nabla f(x^k)^T$$

- Many solutions to choosing step size η^k

1. Fixed step size

- Fixed step size
 - Use fixed value for η^k

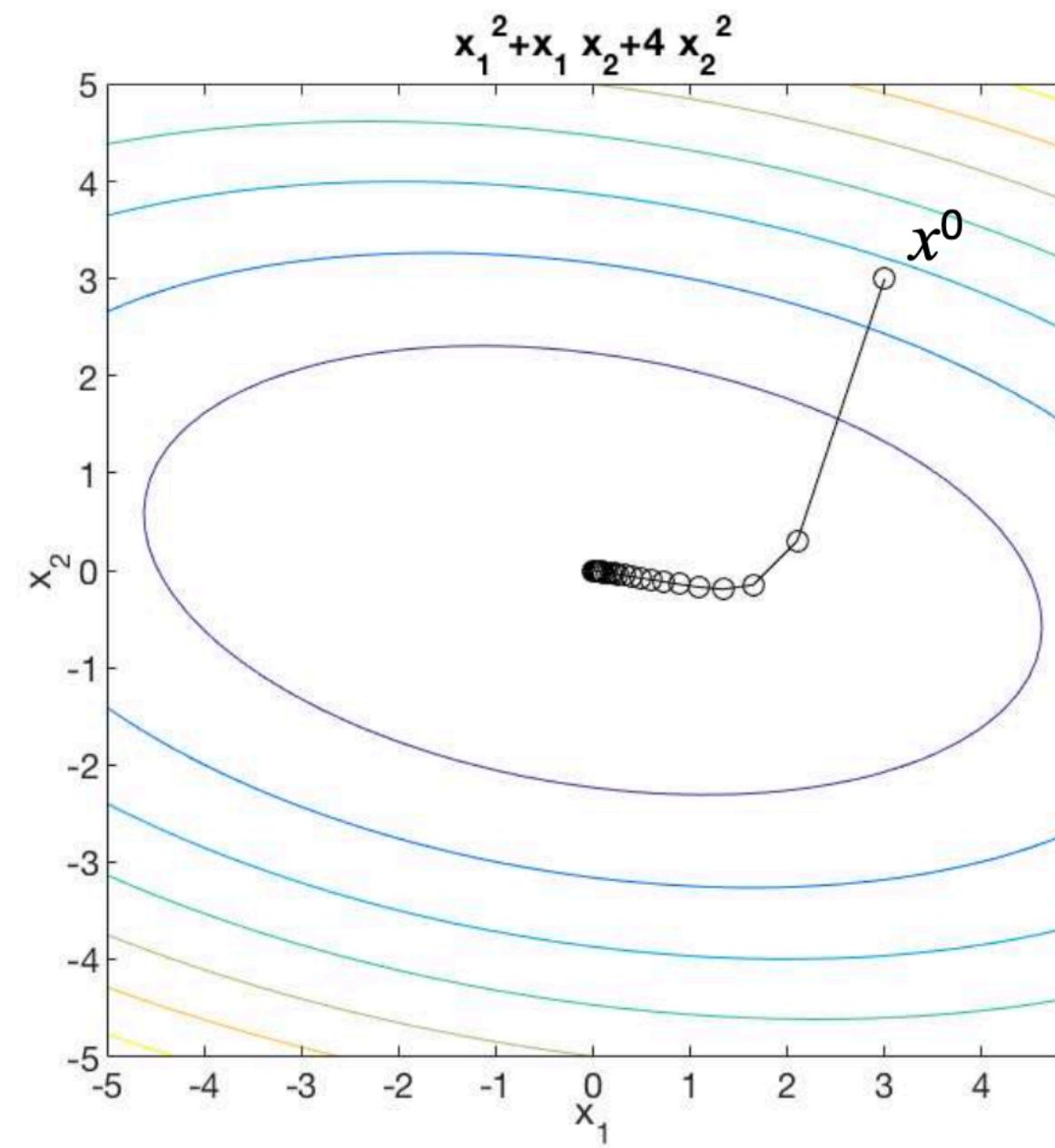


Influence of step size example (constant step size)

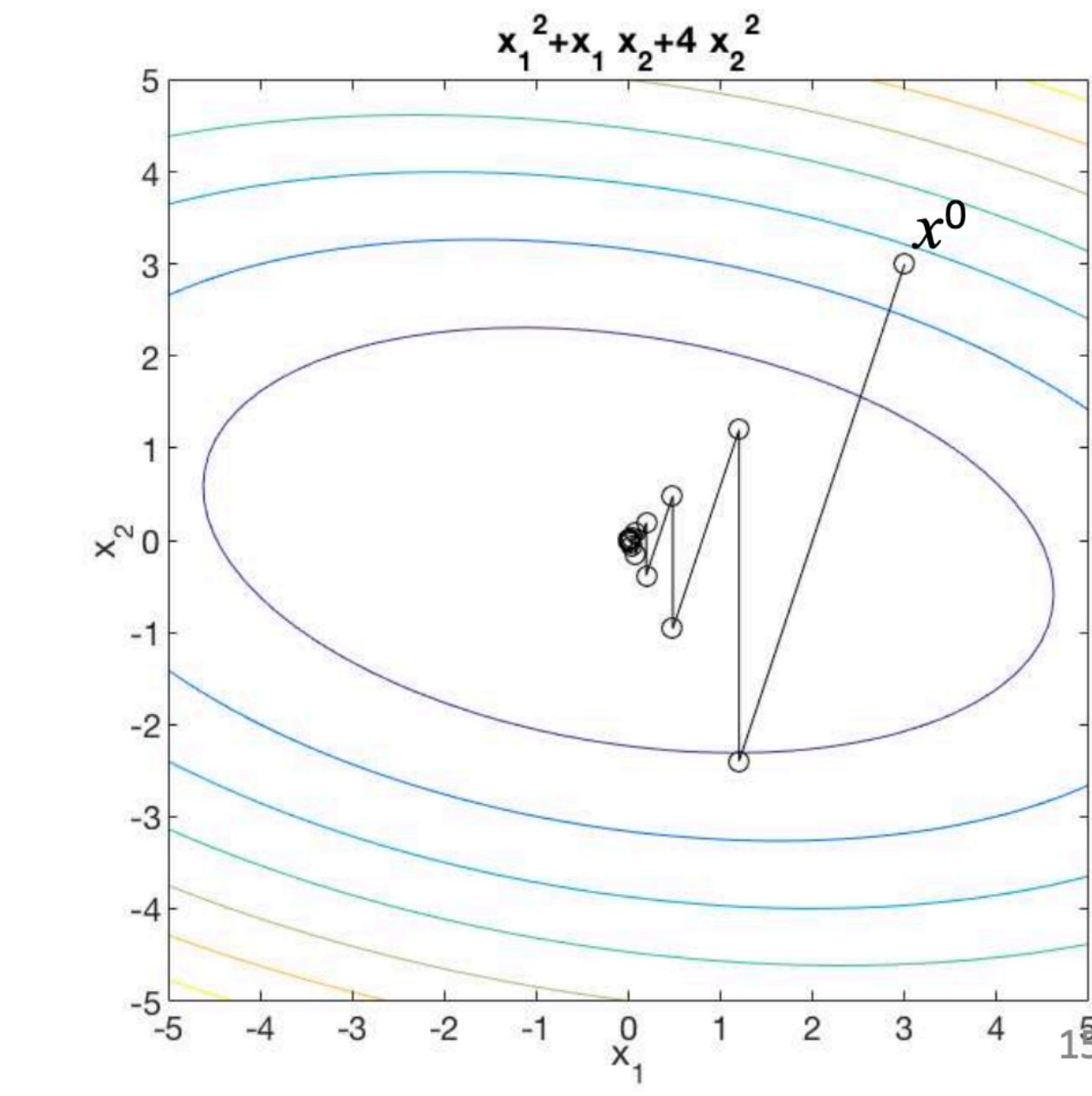
$$f(x_1, x_2) = (x_1)^2 + x_1 x_2 + 4(x_2)^2$$

$$x^{initial} = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

$$\eta = 0.1$$



$$\eta = 0.2$$

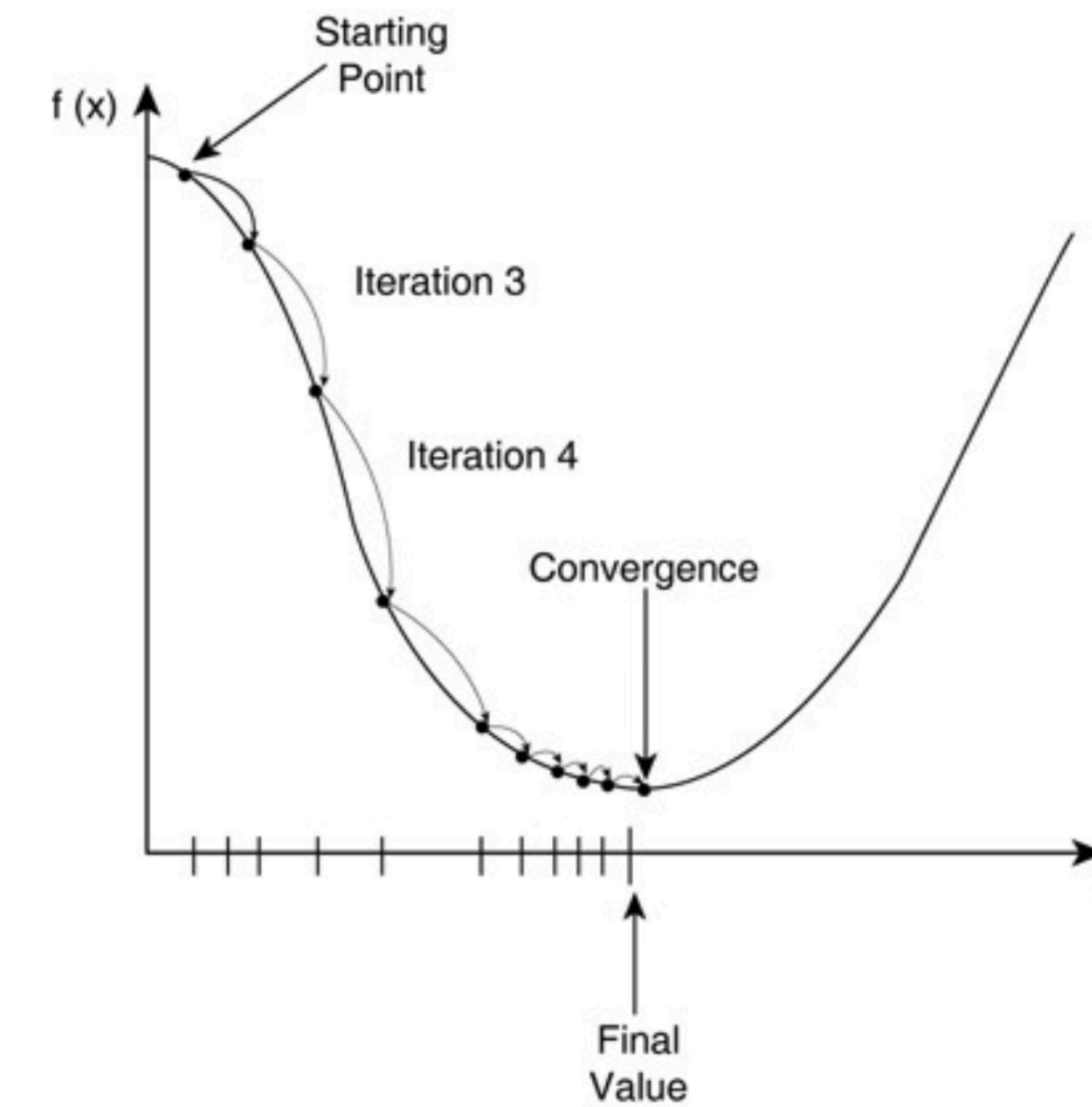


Gradient descent convergence criteria

- The gradient descent algorithm converges when one of the following criteria is satisfied

$$|f(x^{k+1}) - f(x^k)| < \varepsilon_1$$

- Or $\|\nabla f(x^k)\| < \varepsilon_2$



Recap: Gradient Descent Algorithm

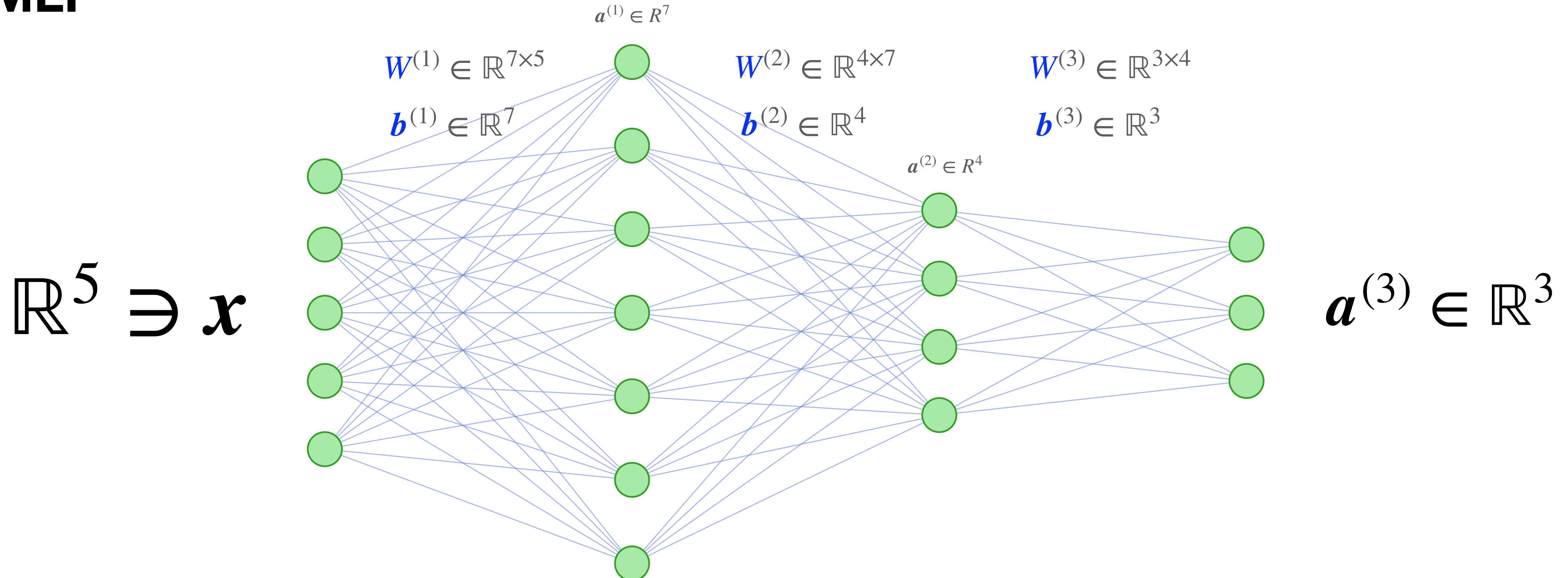
- In order to minimize any function $f(x)$ w.r.t. x
- Initialize:
 - x^0
 - $k = 0$
- While $|f(x^{k+1}) - f(x^k)| > \varepsilon$
 - For every component i
 - $x_i^{k+1} = x_i^k - \eta^k \frac{df}{dx_i}$ Explicitly stating it by component
 - $k = k + 1$

How to calculate gradients?

- Solution = **forward** + **backward** propagation
- Forward propagation -> calculate neuron values (from left to right)
- Backward propagation -> calculate partial derivate of **loss function** w.r.t to neuron values -> calculate partial derivatives w.r.t to **model parameters**.

Forward propagation

MLP



Efficient
implementation

$$x^T \in \mathbb{R}^{1 \times 5}$$

$$z^{(1)} = x^T W^{(1)T} + b^{(1)} \in \mathbb{R}^{1 \times 7}$$

$$a^{(1)} = g_1(z^{(1)}) \in \mathbb{R}^{1 \times 7}$$

$$z^{(2)} = a^{(1)} W^{(2)T} + b^{(2)} \in \mathbb{R}^{1 \times 4}$$

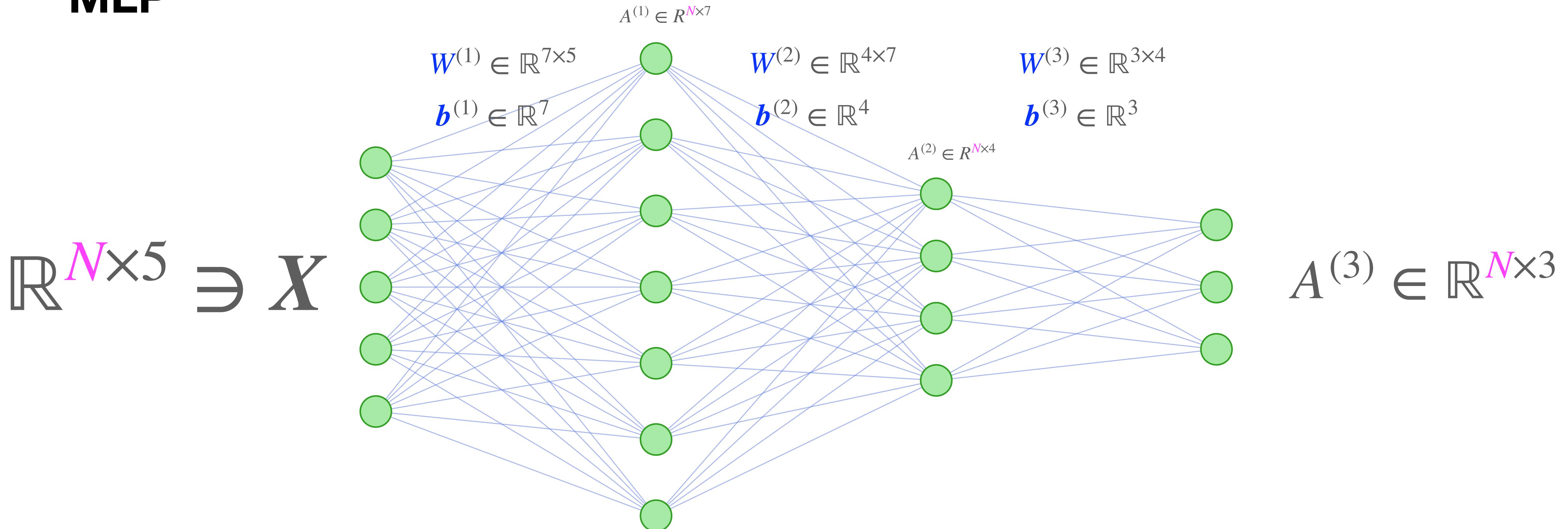
$$a^{(2)} = g_2(z^{(2)}) \in \mathbb{R}^{1 \times 4}$$

$$z^{(3)} = a^{(2)} W^{(3)T} + b^{(3)} \in \mathbb{R}^{1 \times 3}$$

$$a^{(3)} = g_3(z^{(3)}) \in \mathbb{R}^{1 \times 3}$$

Forward propagation

MLP



Efficient
implementation

$$X \in \mathbb{R}^{N \times 5}$$

$$Z^{(1)} = X W^{(1)T} + b^{(1)} \in \mathbb{R}^{N \times 7}$$

$$A^{(1)} = g_1(Z^{(1)}) \in \mathbb{R}^{N \times 7}$$

$$Z^{(2)} = A^{(1)} W^{(2)T} + b^{(2)} \in \mathbb{R}^{N \times 4}$$

$$A^{(2)} = g_2(Z^{(2)}) \in \mathbb{R}^{N \times 4}$$

$$Z^{(3)} = A^{(2)} W^{(3)T} + b^{(3)} \in \mathbb{R}^{N \times 3}$$

$$A^{(3)} = g_3(Z^{(3)}) \in \mathbb{R}^{N \times 3}$$

Let's calculate the derivatives

$$\frac{\partial L}{\partial A^{(p)}} = ?$$

$$L = \frac{1}{N} \sum_{p=1}^N \textcolor{red}{div}(\hat{y}_p, y_p)$$

- Observe that the total loss is the average of the losses of independent samples
- One can do calculations row-wise,
i.e.

$$\frac{\partial L}{\partial A^{(3)}} = -\frac{1}{N} \begin{bmatrix} \frac{\partial L}{\partial a_1^{(3)}} \\ \frac{\partial L}{\partial a_2^{(3)}} \\ \vdots \\ \frac{\partial L}{\partial a_N^{(3)}} \end{bmatrix}, \quad \frac{\partial L}{\partial Z^{(3)}} = -\frac{1}{N} \begin{bmatrix} \frac{\partial L}{\partial z_1^{(3)}} \\ \frac{\partial L}{\partial z_2^{(3)}} \\ \vdots \\ \frac{\partial L}{\partial z_N^{(3)}} \end{bmatrix}, \quad \dots$$

Let's calculate the derivatives

Single sample (row)

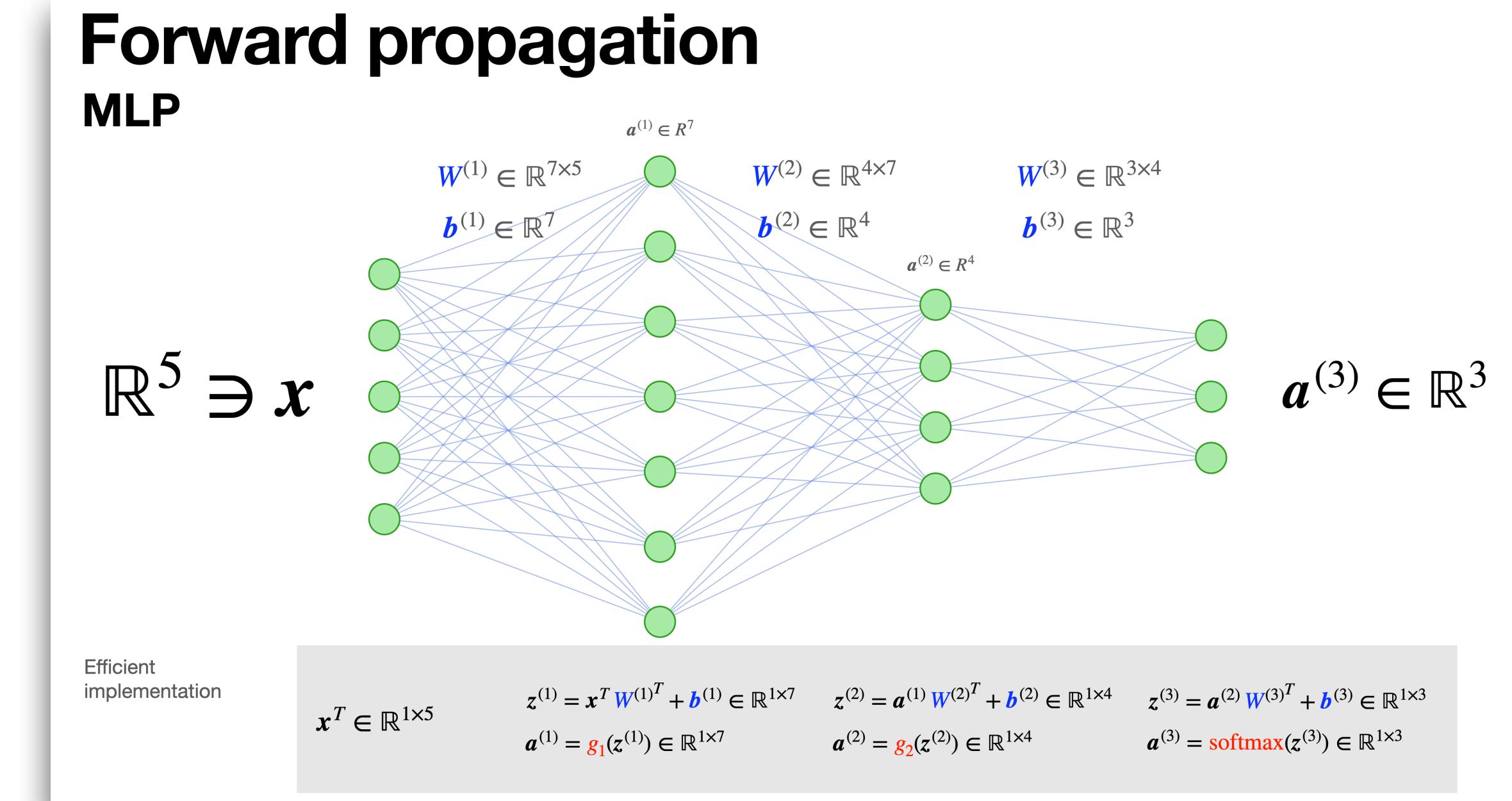
$$L = \frac{1}{N} \sum_{p=1}^N \text{div}(\hat{y}_p, y_p) = \frac{1}{N} \sum_{p=1}^N L_p$$

$$\mathbf{y}^T := Y[1, :]$$

$$\mathbf{a}^{(3)T} := A^{(3)}[1, :] = \mathbf{a}_1^{(3)}$$

$$\mathbf{z}^{(3)T} := Z^{(3)}[1, :]$$

$$\frac{\partial L}{\partial \mathbf{a}^{(3)}} = ?, \quad \frac{\partial L}{\partial \mathbf{z}^{(3)}} = ?, \quad \frac{\partial L}{\partial \mathbf{a}^{(2)}} = ?, \quad \frac{\partial L}{\partial \mathbf{z}^{(2)}} = ?, \quad \frac{\partial L}{\partial \mathbf{a}^{(1)}} = ?, \quad \frac{\partial L}{\partial \mathbf{z}^{(1)}} = ?,$$



Let's calculate the derivatives

Case: MSE loss + “id” activation function

$$\begin{array}{llll} \boldsymbol{x}^T \in \mathbb{R}^{1 \times 5} & z^{(1)} = \boldsymbol{x}^T \mathbf{W}^{(1)T} + \mathbf{b}^{(1)} \in \mathbb{R}^{1 \times 7} & z^{(2)} = \boldsymbol{a}^{(1)} \mathbf{W}^{(2)T} + \mathbf{b}^{(2)} \in \mathbb{R}^{1 \times 4} & z^{(3)} = \boldsymbol{a}^{(2)} \mathbf{W}^{(3)T} + \mathbf{b}^{(3)} \in \mathbb{R}^{1 \times 3} \\ & \boldsymbol{a}^{(1)} = g_1(z^{(1)}) \in \mathbb{R}^{1 \times 7} & \boldsymbol{a}^{(2)} = g_2(z^{(2)}) \in \mathbb{R}^{1 \times 4} & \boldsymbol{a}^{(3)} = \text{id}(z^{(3)}) \in \mathbb{R}^{1 \times 3} \end{array}$$

$$\boldsymbol{y}^T := Y[1, :]$$

$$\boldsymbol{a}^{(3)T} := A^{(3)}[1, :]$$

$$L_1 := \frac{1}{2} \| (Y[1, :] - A^{(3)}[1, :]) \|_2^2$$

$$\bullet \quad \frac{\partial L_1}{\partial \boldsymbol{a}^{(3)}} = ?$$

Let's calculate the derivatives

Case: MSE loss + “id” activation function

$$\begin{array}{llll} \boldsymbol{x}^T \in \mathbb{R}^{1 \times 5} & z^{(1)} = \boldsymbol{x}^T \mathbf{W}^{(1)T} + \mathbf{b}^{(1)} \in \mathbb{R}^{1 \times 7} & z^{(2)} = \boldsymbol{a}^{(1)} \mathbf{W}^{(2)T} + \mathbf{b}^{(2)} \in \mathbb{R}^{1 \times 4} & z^{(3)} = \boldsymbol{a}^{(2)} \mathbf{W}^{(3)T} + \mathbf{b}^{(3)} \in \mathbb{R}^{1 \times 3} \\ & \boldsymbol{a}^{(1)} = \mathbf{g}_1(z^{(1)}) \in \mathbb{R}^{1 \times 7} & \boldsymbol{a}^{(2)} = \mathbf{g}_2(z^{(2)}) \in \mathbb{R}^{1 \times 4} & \boldsymbol{a}^{(3)} = \mathbf{id}(z^{(3)}) \in \mathbb{R}^{1 \times 3} \end{array}$$

$$\boldsymbol{y}^T := Y[1, :]$$

$$\boldsymbol{a}^{(3)T} := A^{(3)}[1, :]$$

$$L_1 := \frac{1}{2} \| (Y[1, :] - A^{(3)}[1, :]) \|_2^2$$

$$\cdot \frac{\partial L_1}{\partial \boldsymbol{a}^{(3)}} = \boldsymbol{a}^{(3)T} - \boldsymbol{y}^T$$

$$\frac{\partial L_1}{\partial z^{(3)}} = \frac{\partial L_1}{\partial \boldsymbol{a}^{(3)}} \cdot \frac{\partial \boldsymbol{a}^{(3)}}{\partial z^{(3)}} =$$

Let's calculate the derivatives

Case: MSE loss + “id” activation function

$$\begin{array}{llll} \boldsymbol{x}^T \in \mathbb{R}^{1 \times 5} & z^{(1)} = \boldsymbol{x}^T \mathbf{W}^{(1)T} + \mathbf{b}^{(1)} \in \mathbb{R}^{1 \times 7} & z^{(2)} = \boldsymbol{a}^{(1)} \mathbf{W}^{(2)T} + \mathbf{b}^{(2)} \in \mathbb{R}^{1 \times 4} & z^{(3)} = \boldsymbol{a}^{(2)} \mathbf{W}^{(3)T} + \mathbf{b}^{(3)} \in \mathbb{R}^{1 \times 3} \\ & \boldsymbol{a}^{(1)} = g_1(z^{(1)}) \in \mathbb{R}^{1 \times 7} & \boldsymbol{a}^{(2)} = g_2(z^{(2)}) \in \mathbb{R}^{1 \times 4} & \boldsymbol{a}^{(3)} = \text{id}(z^{(3)}) \in \mathbb{R}^{1 \times 3} \end{array}$$

$$\boldsymbol{y}^T := Y[1, :]$$

$$\boldsymbol{a}^{(3)T} := A^{(3)}[1, :]$$

$$L_1 := \frac{1}{2} \| (Y[1, :] - A^{(3)}[1, :]) \|_2^2$$

$$\cdot \frac{\partial L_1}{\partial \boldsymbol{a}^{(3)}} = \boldsymbol{a}^{(3)T} - \boldsymbol{y}^T$$

$$\frac{\partial L_1}{\partial z^{(3)}} = \frac{\partial L_1}{\partial \boldsymbol{a}^{(3)}} \cdot \frac{\partial \boldsymbol{a}^{(3)}}{\partial z^{(3)}} = \boldsymbol{a}^{(3)T} - \boldsymbol{y}^T$$

Let's calculate the derivatives

Case: MSE loss + “id” activation function

$$\begin{array}{llll} \boldsymbol{x}^T \in \mathbb{R}^{1 \times 5} & z^{(1)} = \boldsymbol{x}^T \mathbf{W}^{(1)T} + \mathbf{b}^{(1)} \in \mathbb{R}^{1 \times 7} & z^{(2)} = \boldsymbol{a}^{(1)} \mathbf{W}^{(2)T} + \mathbf{b}^{(2)} \in \mathbb{R}^{1 \times 4} & z^{(3)} = \boldsymbol{a}^{(2)} \mathbf{W}^{(3)T} + \mathbf{b}^{(3)} \in \mathbb{R}^{1 \times 3} \\ & \boldsymbol{a}^{(1)} = \mathbf{g}_1(z^{(1)}) \in \mathbb{R}^{1 \times 7} & \boldsymbol{a}^{(2)} = \mathbf{g}_2(z^{(2)}) \in \mathbb{R}^{1 \times 4} & \boldsymbol{a}^{(3)} = \mathbf{id}(z^{(3)}) \in \mathbb{R}^{1 \times 3} \end{array}$$

$$\boldsymbol{y}^T := Y[1, :]$$

$$\boldsymbol{a}^{(3)T} := A^{(3)}[1, :]$$

$$L_1 := \frac{1}{2} \| (Y[1, :] - A^{(3)}[1, :]) \|_2^2$$

$$\cdot \frac{\partial L_1}{\partial \boldsymbol{a}^{(3)}} = \boldsymbol{a}^{(3)T} - \boldsymbol{y}^T$$

$$\frac{\partial L_1}{\partial z^{(3)}} = \frac{\partial L_1}{\partial \boldsymbol{a}^{(3)}} \cdot \frac{\partial \boldsymbol{a}^{(3)}}{\partial z^{(3)}} = \boldsymbol{a}^{(3)T} - \boldsymbol{y}^T$$

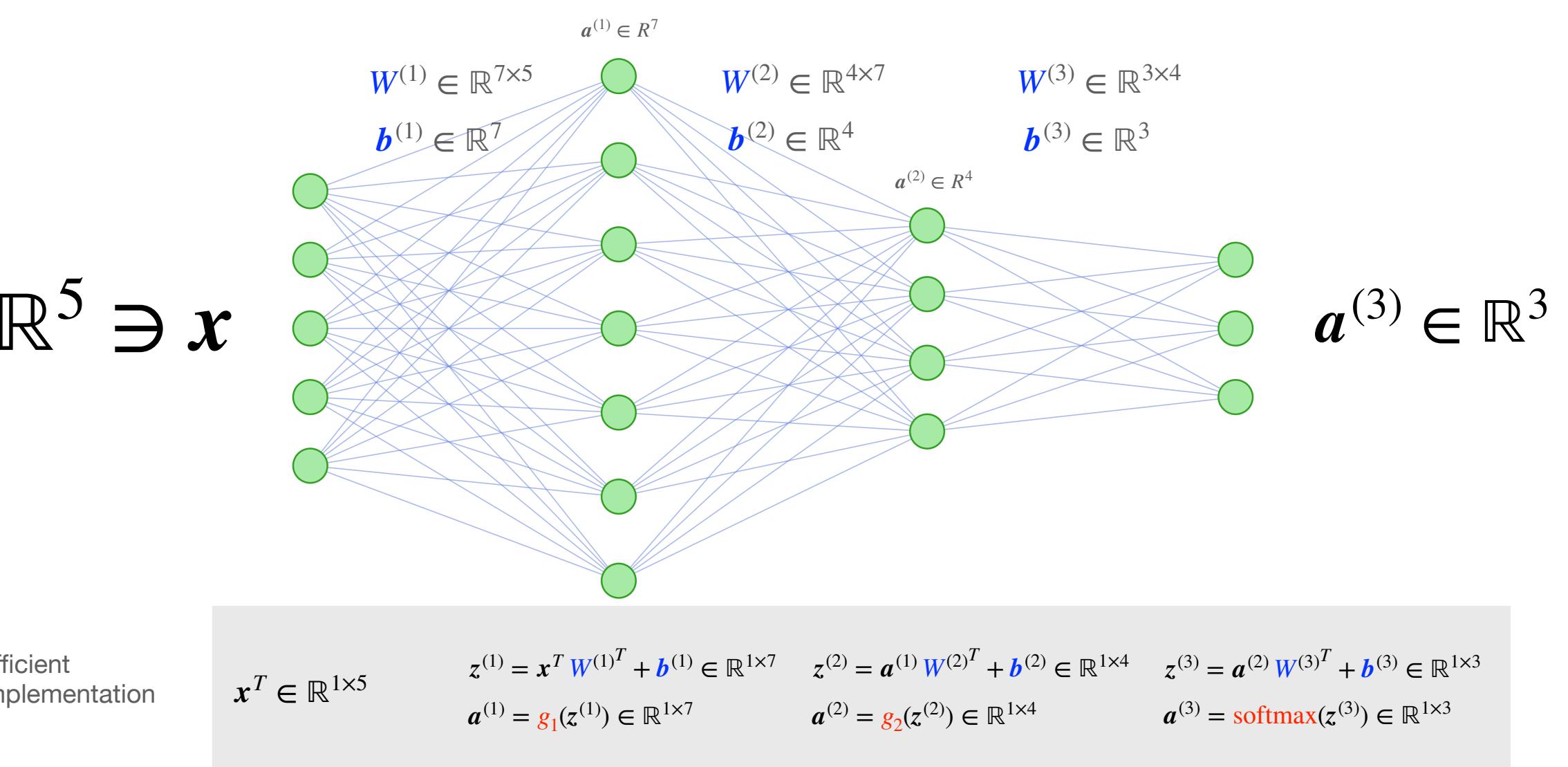
Matrix case ->

$$\frac{\partial L}{\partial Z^{(3)}} = \frac{1}{N} (A^{(3)} - Y)$$

Backward propagation

MLP

$$L_1 := \frac{1}{2} \| (Y[1, :] - A^{(3)}[1, :]) \|_2^2$$



$$\frac{\partial L_1}{\partial a^{(3)}} = ???$$

$$\frac{\partial L_1}{\partial z^{(3)}} = a^{(3)T} - y^T \in \mathbb{R}^{1 \times 3}$$

$$\frac{\partial L_1}{\partial W^{(3)}} = \frac{\partial L_1}{\partial z^{(3)}}^T \cdot \frac{\partial z^{(3)T}}{\partial W^{(3)}} = \frac{\partial L_1}{\partial z^{(3)}}^T \cdot a^{(2)T}$$

$$\frac{\partial L_1}{\partial a^{(2)}} = \frac{\partial L_1}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} = \frac{\partial L_1}{\partial z^{(3)}} \cdot W^{(3)}$$

$$\frac{\partial L_1}{\partial z^{(2)}} = \frac{\partial L_1}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} = \frac{\partial L_1}{\partial a^{(2)}} \otimes \frac{\partial g_2}{\partial z^{(2)}}$$

$$\frac{\partial L_1}{\partial W^{(2)}} = \frac{\partial L_1}{\partial z^{(2)}}^T \cdot \frac{\partial z^{(2)}}{\partial W^{(2)}} = \frac{\partial L_1}{\partial z^{(2)}}^T \cdot a^{(1)T}$$

$$\frac{\partial L_1}{\partial a^{(1)}} = \frac{\partial L_1}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial a^{(1)}} = \frac{\partial L_1}{\partial z^{(2)}} \cdot W^{(2)}$$

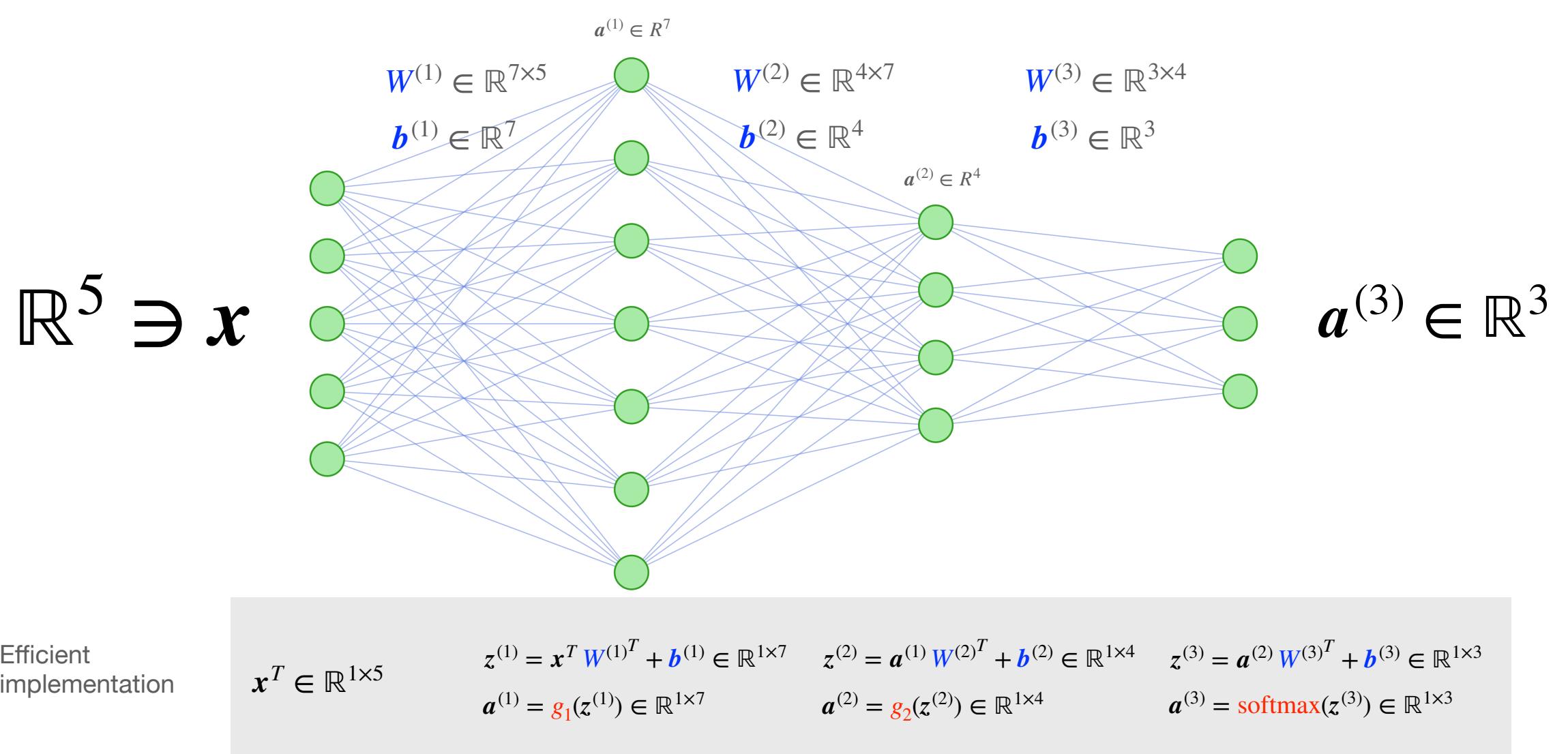
$$\frac{\partial L_1}{\partial z^{(1)}} = \frac{\partial L_1}{\partial a^{(1)}} \cdot \frac{\partial a^{(1)}}{\partial z^{(1)}} = \frac{\partial L_1}{\partial a^{(1)}} \otimes \frac{\partial g_1}{\partial z^{(1)}}$$

$$\frac{\partial L_1}{\partial W^{(1)}} = \frac{\partial L_1}{\partial z^{(1)}}^T \cdot \frac{\partial z^{(1)}}{\partial W^{(1)}} = \frac{\partial L_1}{\partial z^{(1)}}^T \cdot x^T$$

Backward propagation

MLP

$$L_1 := \frac{1}{2} \| (Y[1, :] - A^{(3)}[1, :]) \|_2^2$$



$$\frac{\partial L_1}{\partial W^{(3)}} = \frac{\partial L_1}{\partial z^{(3)}}^T \cdot \frac{\partial z^{(3)T}}{\partial W^{(3)}} = \frac{\partial L_1}{\partial z^{(3)}}^T \cdot a^{(2)T}$$

$$\frac{\partial L_1}{\partial b^{(3)}} = \frac{\partial L_1}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial b^{(3)}} = \frac{\partial L_1}{\partial z^{(3)}}$$

$$\frac{\partial L_1}{\partial W^{(2)}} = \frac{\partial L_1}{\partial z^{(2)}}^T \cdot \frac{\partial z^{(2)}}{\partial W^{(2)}} = \frac{\partial L_1}{\partial z^{(2)}}^T \cdot a^{(1)T}$$

$$\frac{\partial L_1}{\partial b^{(2)}} = \frac{\partial L_1}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial b^{(2)}} = \frac{\partial L_1}{\partial z^{(2)}}$$

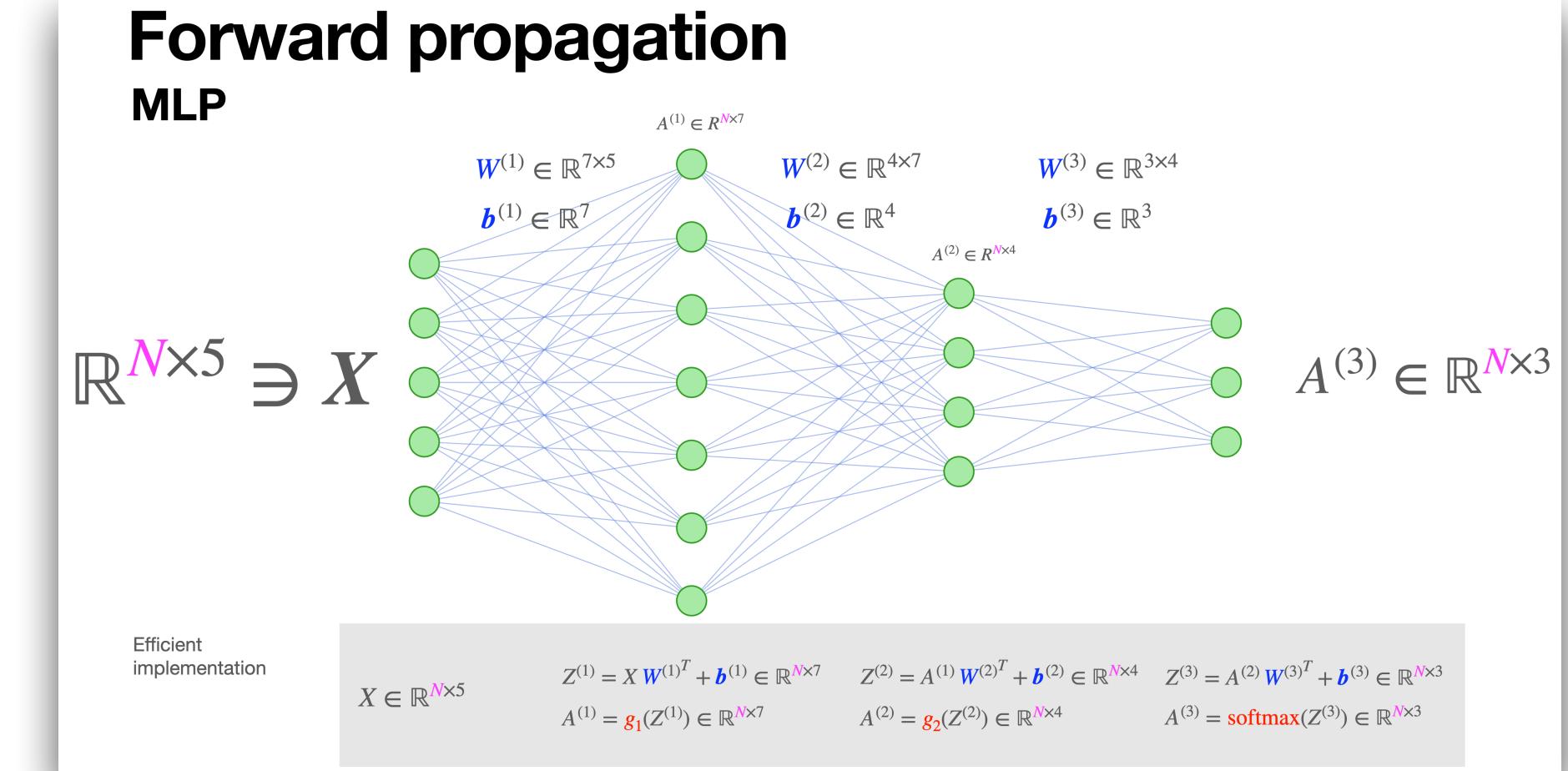
$$\frac{\partial L_1}{\partial W^{(1)}} = \frac{\partial L_1}{\partial z^{(1)}}^T \cdot \frac{\partial z^{(1)}}{\partial W^{(1)}} = \frac{\partial L_1}{\partial z^{(1)}}^T \cdot x^T$$

$$\frac{\partial L_1}{\partial b^{(1)}} = \frac{\partial L_1}{\partial z^{(1)}} \cdot \frac{\partial z^{(1)}}{\partial b^{(1)}} = \frac{\partial L_1}{\partial z^{(1)}}$$

Backward propagation

MLP

$$L = -\frac{1}{N} \sum_{p=1}^N L_p$$



$$\frac{\partial L}{\partial W^{(3)}} = \frac{\partial L}{\partial Z^{(3)}}^T \cdot \frac{\partial Z^{(3)}}{\partial W^{(3)}} = \frac{\partial L}{\partial Z^{(3)}}^T \cdot Z^{(3)} \cdot \frac{\partial}{\partial W^{(3)}} = \frac{\partial L}{\partial Z^{(3)}}^T \cdot A^{(2)}$$

$$\frac{\partial L_1}{\partial b^{(3)}} = \frac{\partial L_1}{\partial Z^{(3)}}^T \cdot \frac{\partial Z^{(3)}}{\partial b^{(3)}} = \frac{\partial L_1}{\partial Z^{(3)}}^T \cdot \mathbf{1}^N$$

$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial Z^{(2)}}^T \cdot \frac{\partial Z^{(2)}}{\partial W^{(2)}} = \frac{\partial L}{\partial Z^{(2)}}^T \cdot Z^{(2)} \cdot \frac{\partial}{\partial W^{(2)}} = \frac{\partial L}{\partial Z^{(2)}}^T \cdot A^{(1)}$$

$$\frac{\partial L_1}{\partial b^{(2)}} = \frac{\partial L_1}{\partial Z^{(2)}}^T \cdot \frac{\partial Z^{(2)}}{\partial b^{(2)}} = \frac{\partial L_1}{\partial Z^{(2)}}^T \cdot \mathbf{1}^N$$

$$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial Z^{(1)}}^T \cdot \frac{\partial Z^{(1)}}{\partial W^{(1)}} = \frac{\partial L}{\partial Z^{(1)}}^T \cdot Z^{(1)} \cdot \frac{\partial}{\partial W^{(1)}} = \frac{\partial L}{\partial Z^{(1)}}^T \cdot X$$

$$\frac{\partial L_1}{\partial b^{(1)}} = \frac{\partial L_1}{\partial Z^{(1)}}^T \cdot \frac{\partial Z^{(1)}}{\partial b^{(1)}} = \frac{\partial L_1}{\partial Z^{(1)}}^T \cdot \mathbf{1}^N$$

Let's calculate the derivatives

Case: cross entropy loss / “sigmoid” activation function

$$\boldsymbol{x}^T \in \mathbb{R}^{1 \times 5}$$

$$z^{(1)} = \boldsymbol{x}^T \mathbf{W}^{(1)T} + \mathbf{b}^{(1)} \in \mathbb{R}^{1 \times 7}$$
$$\boldsymbol{a}^{(1)} = \mathbf{g}_1(z^{(1)}) \in \mathbb{R}^{1 \times 7}$$

$$z^{(2)} = \boldsymbol{a}^{(1)} \mathbf{W}^{(2)T} + \mathbf{b}^{(2)} \in \mathbb{R}^{1 \times 4}$$
$$\boldsymbol{a}^{(2)} = \mathbf{g}_2(z^{(2)}) \in \mathbb{R}^{1 \times 4}$$

$$z^{(3)} = \boldsymbol{a}^{(2)} \mathbf{W}^{(3)T} + \mathbf{b}^{(3)} \in \mathbb{R}^{1 \times 3}$$
$$\boldsymbol{a}^{(3)} = \mathbf{o}(z^{(3)}) \in \mathbb{R}^{1 \times 3}$$

Logistic regression

Derivative of loss function

$$z = \mathbf{x}\mathbf{w}^T + b, \quad \hat{y} = \sigma(z), \quad \frac{\partial \hat{y}}{\partial z} = \sigma(z)(1 - \sigma'(z))$$

$$L(\mathbf{w}, b) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

$$\frac{\partial L}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1 - y}{1 - \hat{y}}, \quad \frac{\partial L}{\partial z} = \hat{y} - y, \quad \frac{\partial L}{\partial \mathbf{w}} = (\hat{y} - y) \mathbf{x}^T, \quad \frac{\partial L}{\partial b} = \hat{y} - y$$