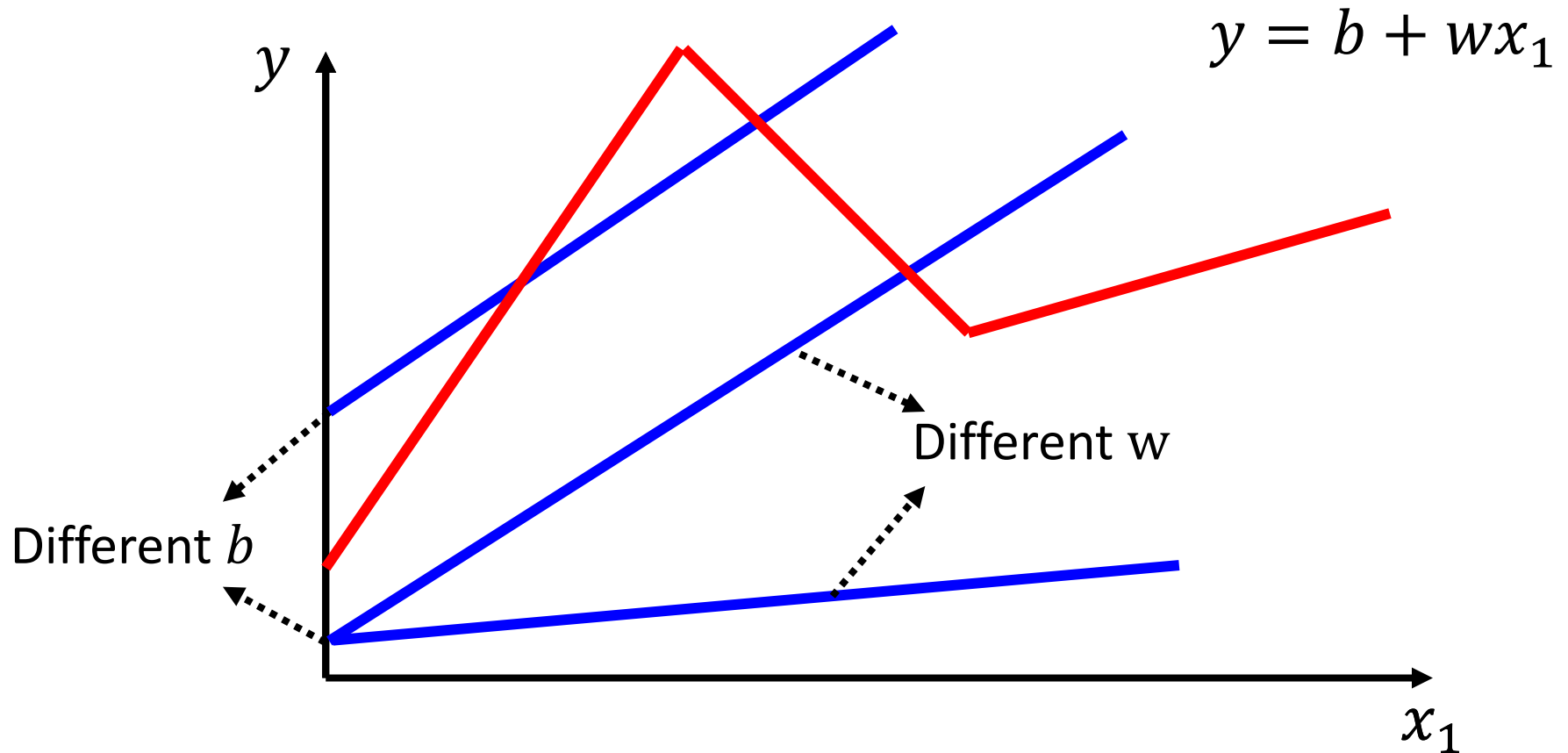


Something before Real Deep  
Learning

- From Linear Model  
to **Neural Network**

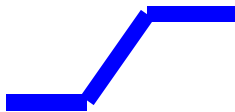
Yizhen Lao

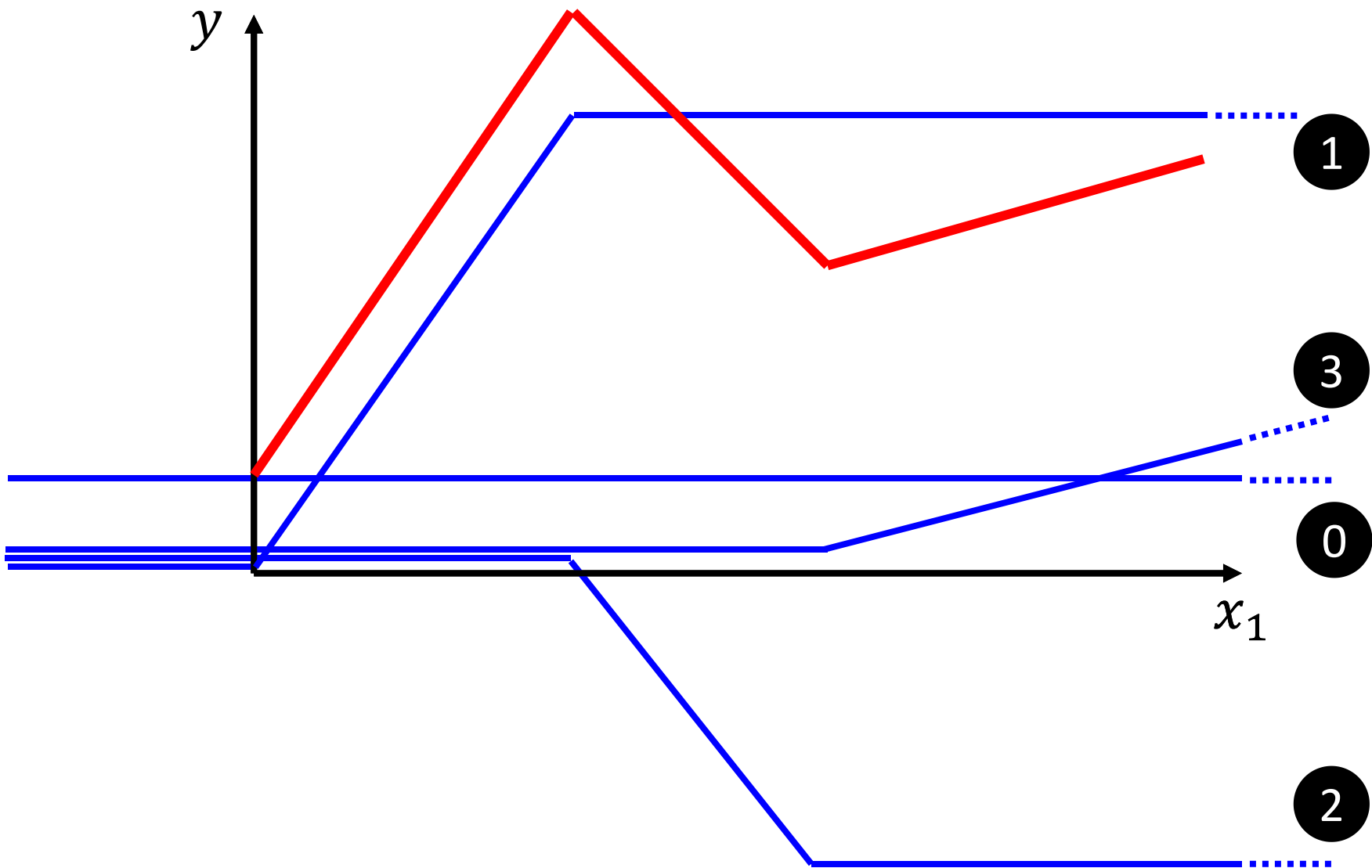
Linear models are too simple ... we need more sophisticated modes.



Linear models have severe limitation. ***Model Bias***

We need a more flexible model!

red curve = constant + sum of a set of 



# All Piecewise Linear Curves

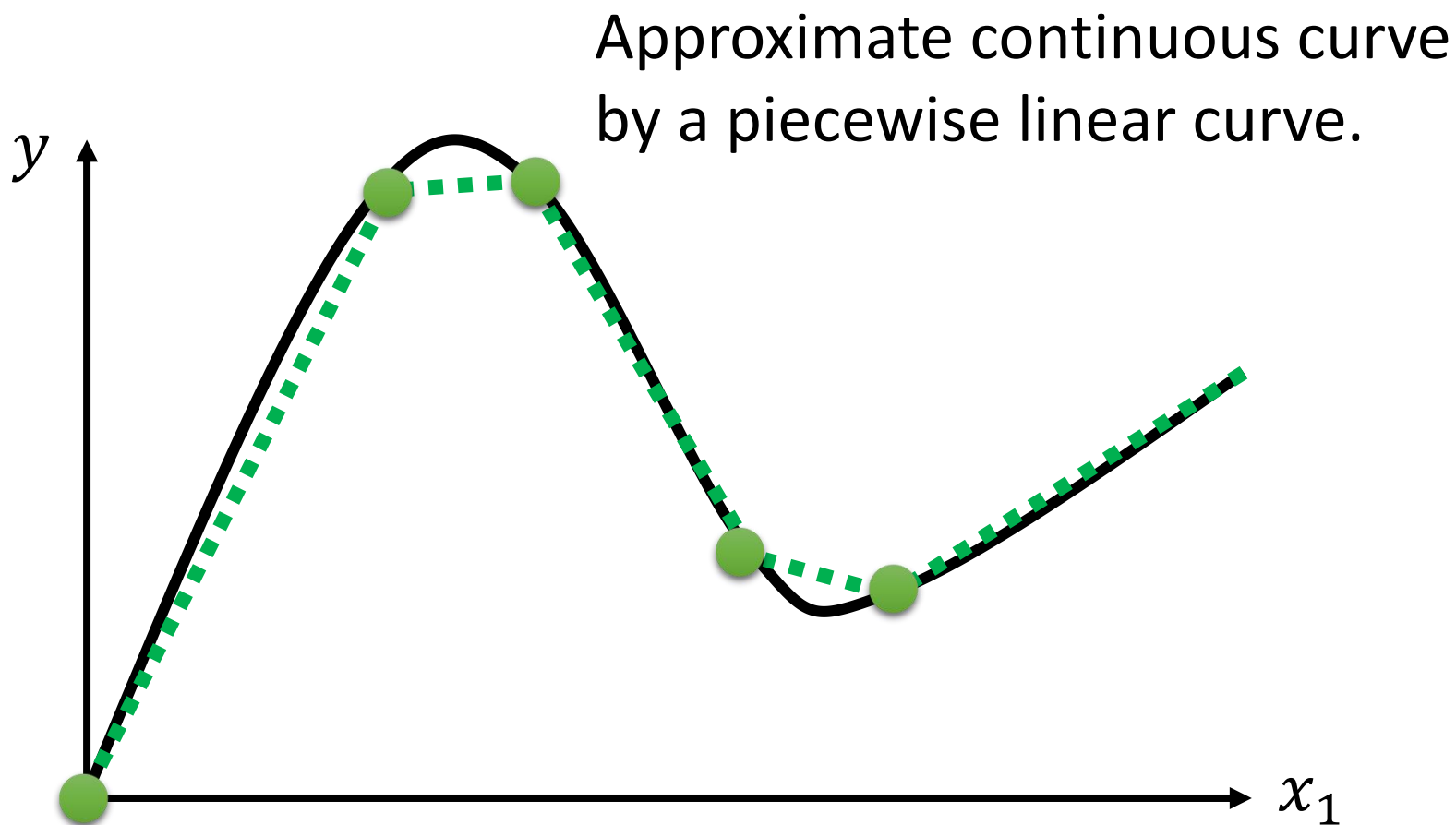
= constant + sum of a set of



More pieces require more

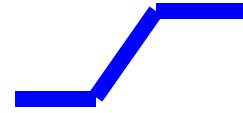


# Beyond Piecewise Linear?



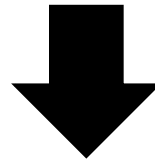
To have good approximation, we need sufficient pieces.

red curve = constant + sum of a set of



How to represent  
this function?

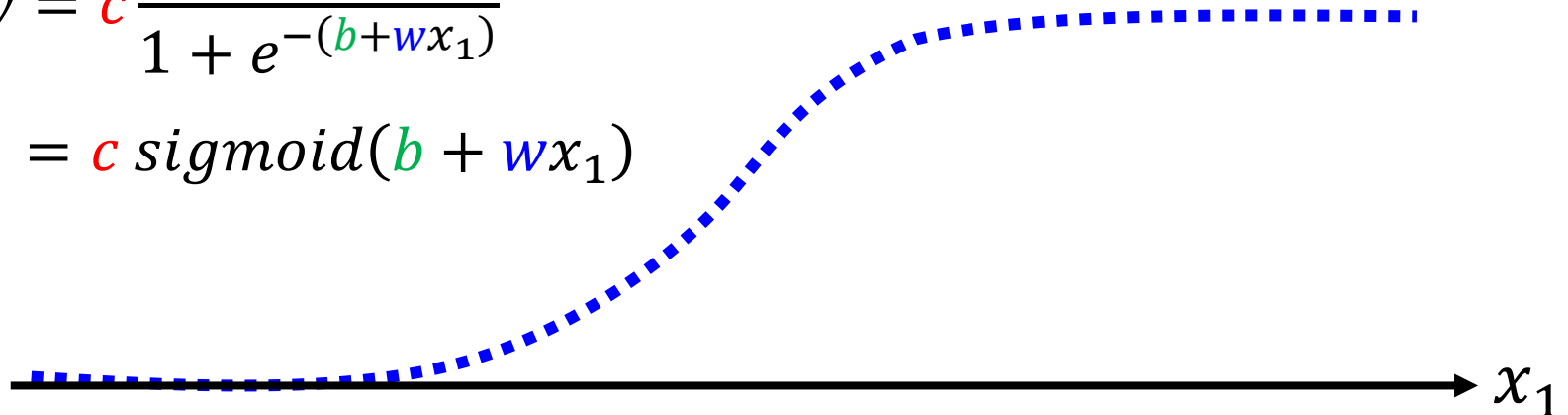
Hard Sigmoid

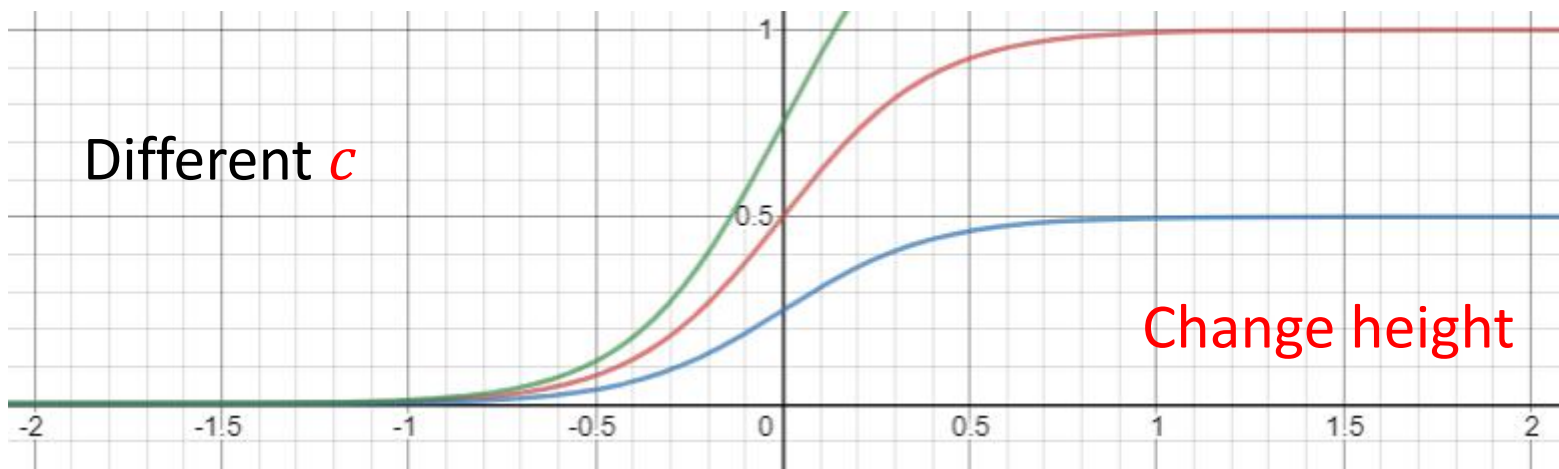
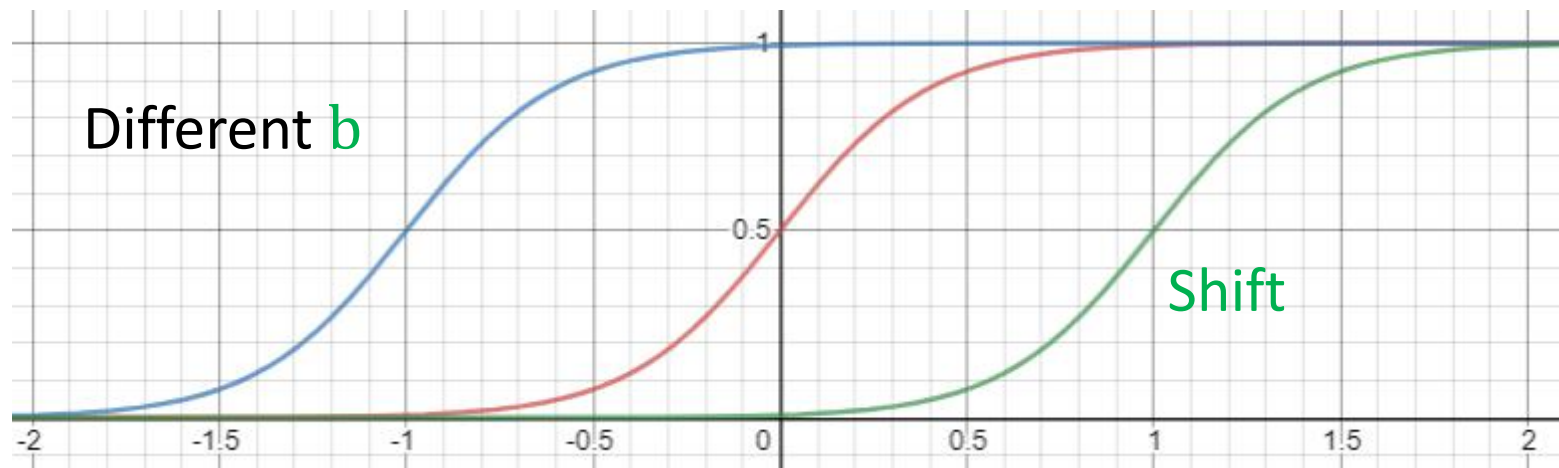
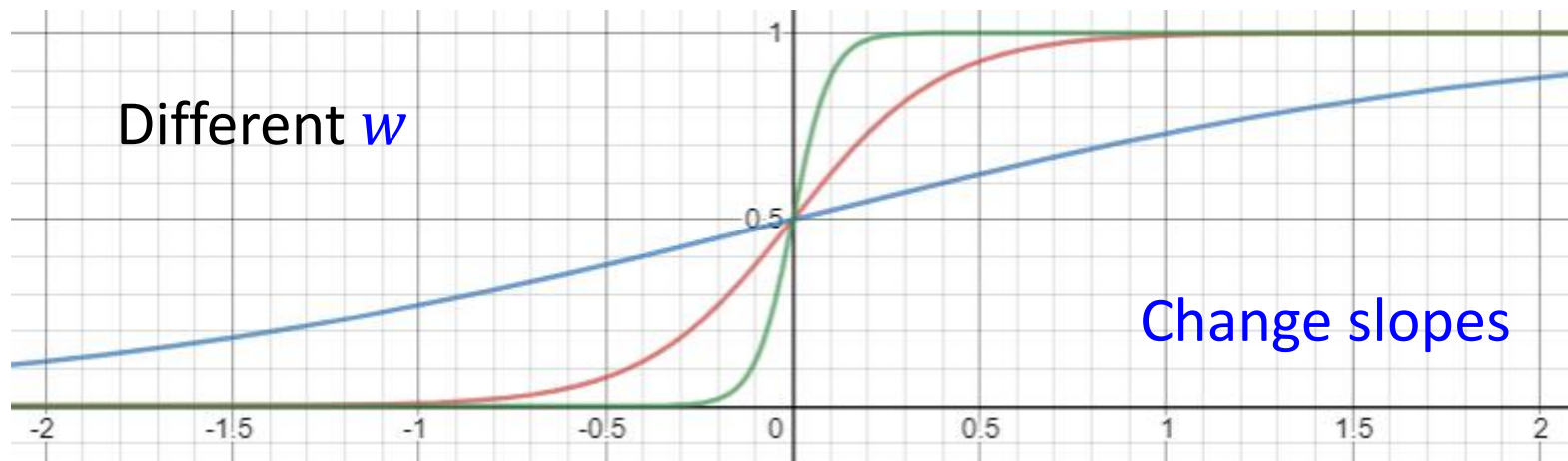


Sigmoid Function

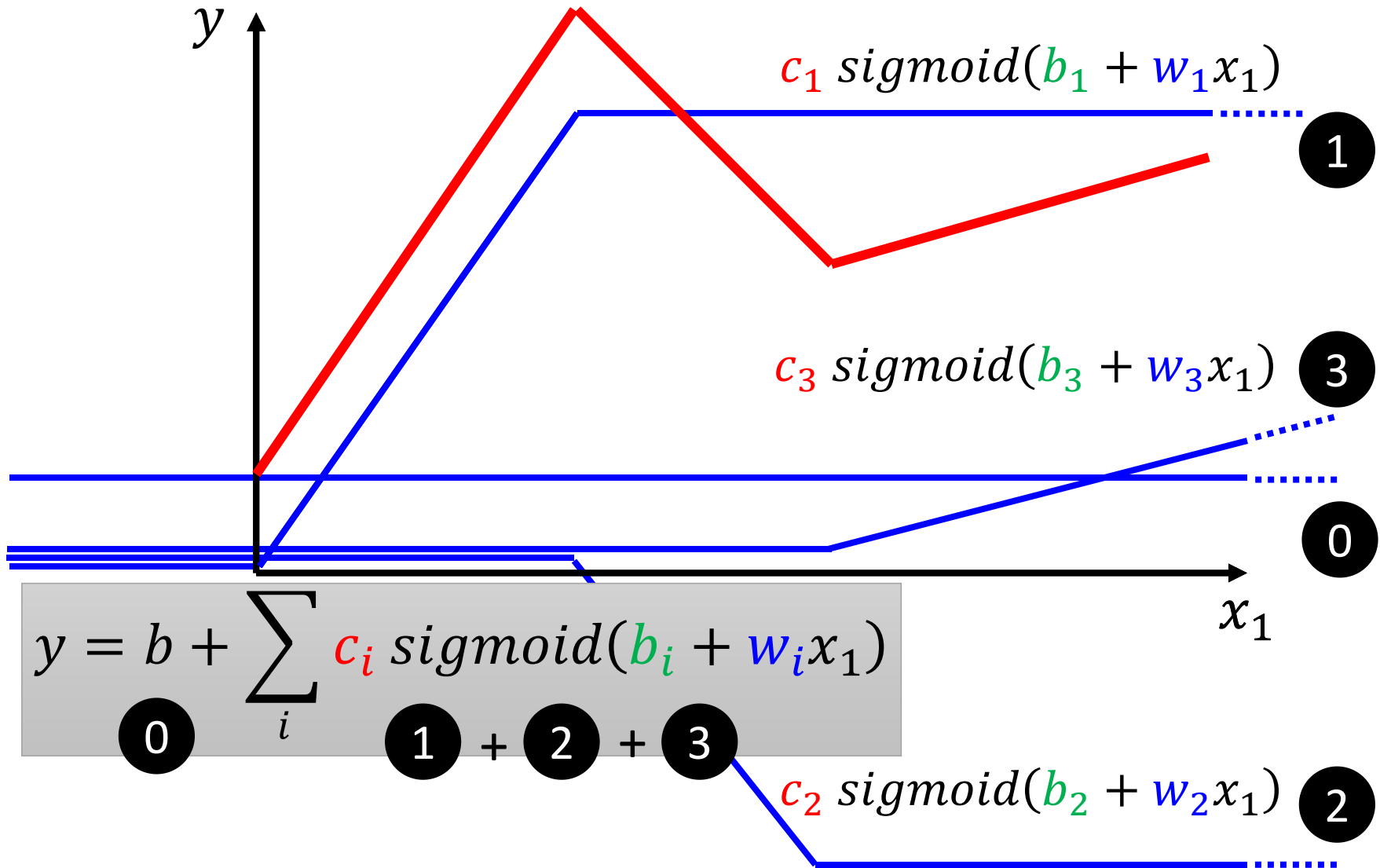
$$y = c \frac{1}{1 + e^{-(b + wx_1)}}$$

$$= c \operatorname{sigmoid}(b + wx_1)$$





red curve = sum of a set of  + constant



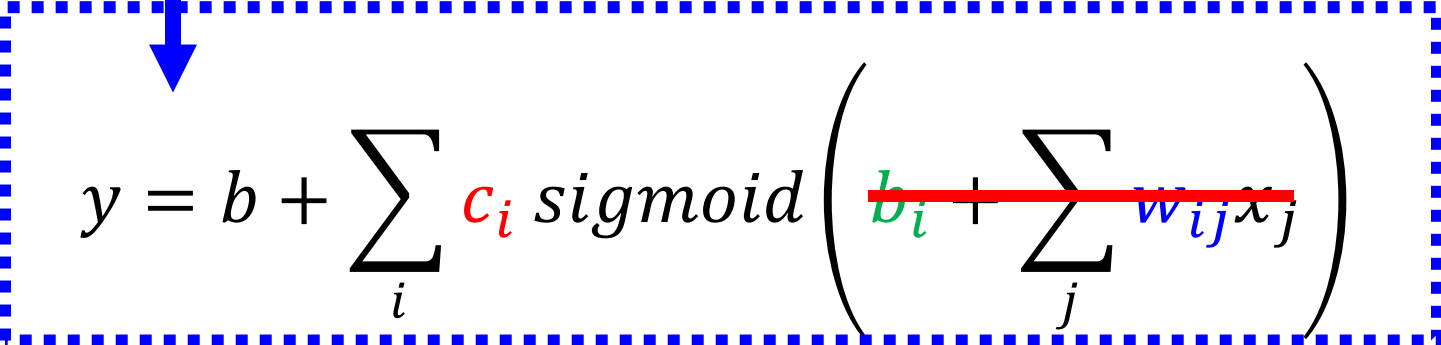



# New Model: More Features

$$y = \underline{b + wx_1}$$

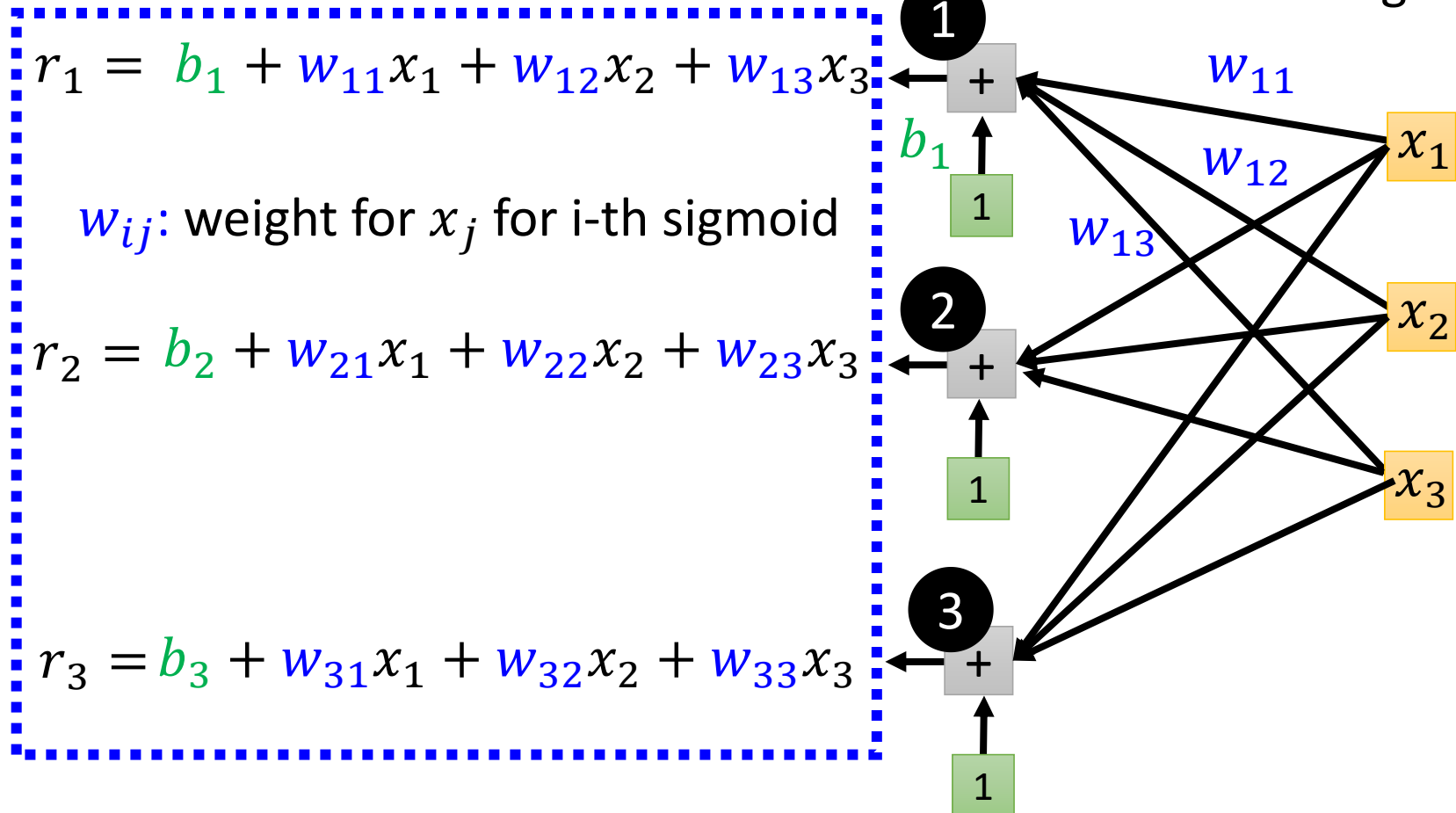

$$y = b + \sum_i c_i \operatorname{sigmoid}(\underline{b_i + w_i x_1})$$

$$y = \underline{b + \sum_j w_j x_j}$$


$$y = b + \sum_i c_i \operatorname{sigmoid}\left(\cancel{b_i} + \sum_j \cancel{w_{ij} x_j}\right)$$

$$y = b + \sum_i \textcolor{red}{c}_i \textit{sigmoid} \left( \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right)$$

$j: 1, 2, 3$   
 no. of features  
 $i: 1, 2, 3$   
 no. of sigmoid



$$y = b + \sum_i \textcolor{red}{c}_i \textit{sigmoid} \left( \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right) \quad \begin{array}{l} i: 1,2,3 \\ j: 1,2,3 \end{array}$$

$$r_1 = \textcolor{green}{b}_1 + \textcolor{blue}{w}_{11}x_1 + \textcolor{blue}{w}_{12}x_2 + \textcolor{blue}{w}_{13}x_3$$

$$r_2 = \textcolor{green}{b}_2 + \textcolor{blue}{w}_{21}x_1 + \textcolor{blue}{w}_{22}x_2 + \textcolor{blue}{w}_{23}x_3$$

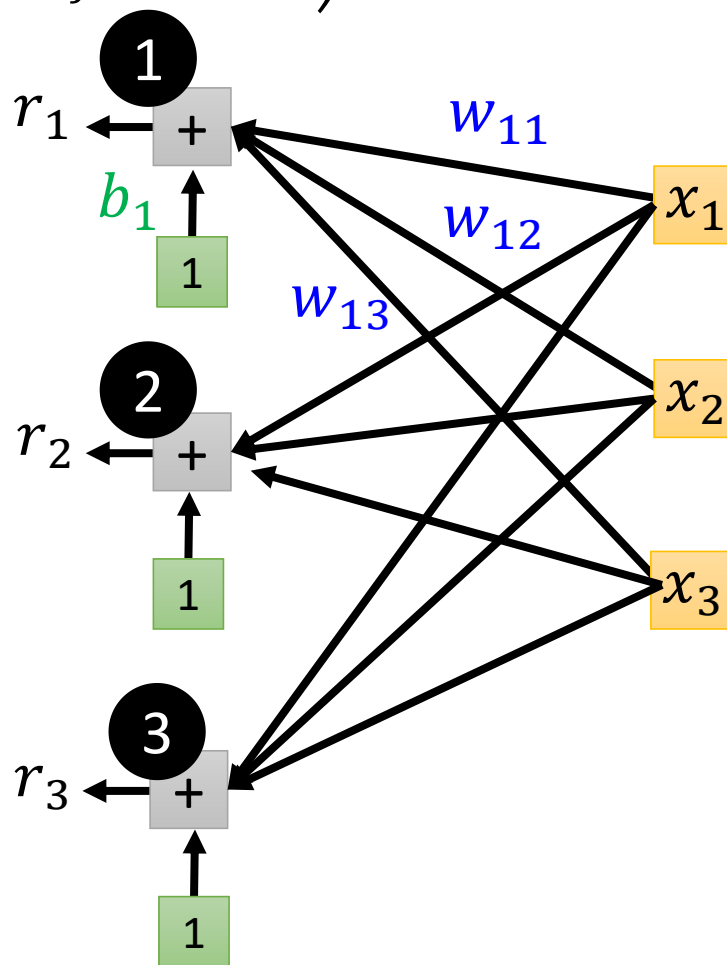
$$r_3 = \textcolor{green}{b}_3 + \textcolor{blue}{w}_{31}x_1 + \textcolor{blue}{w}_{32}x_2 + \textcolor{blue}{w}_{33}x_3$$

$$\begin{bmatrix} r_1 \\ r_2 \\ r_3 \end{bmatrix} = \begin{bmatrix} \textcolor{green}{b}_1 \\ \textcolor{green}{b}_2 \\ \textcolor{green}{b}_3 \end{bmatrix} + \begin{bmatrix} \textcolor{blue}{w}_{11} & \textcolor{blue}{w}_{12} & \textcolor{blue}{w}_{13} \\ \textcolor{blue}{w}_{21} & \textcolor{blue}{w}_{22} & \textcolor{blue}{w}_{23} \\ \textcolor{blue}{w}_{31} & \textcolor{blue}{w}_{32} & \textcolor{blue}{w}_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

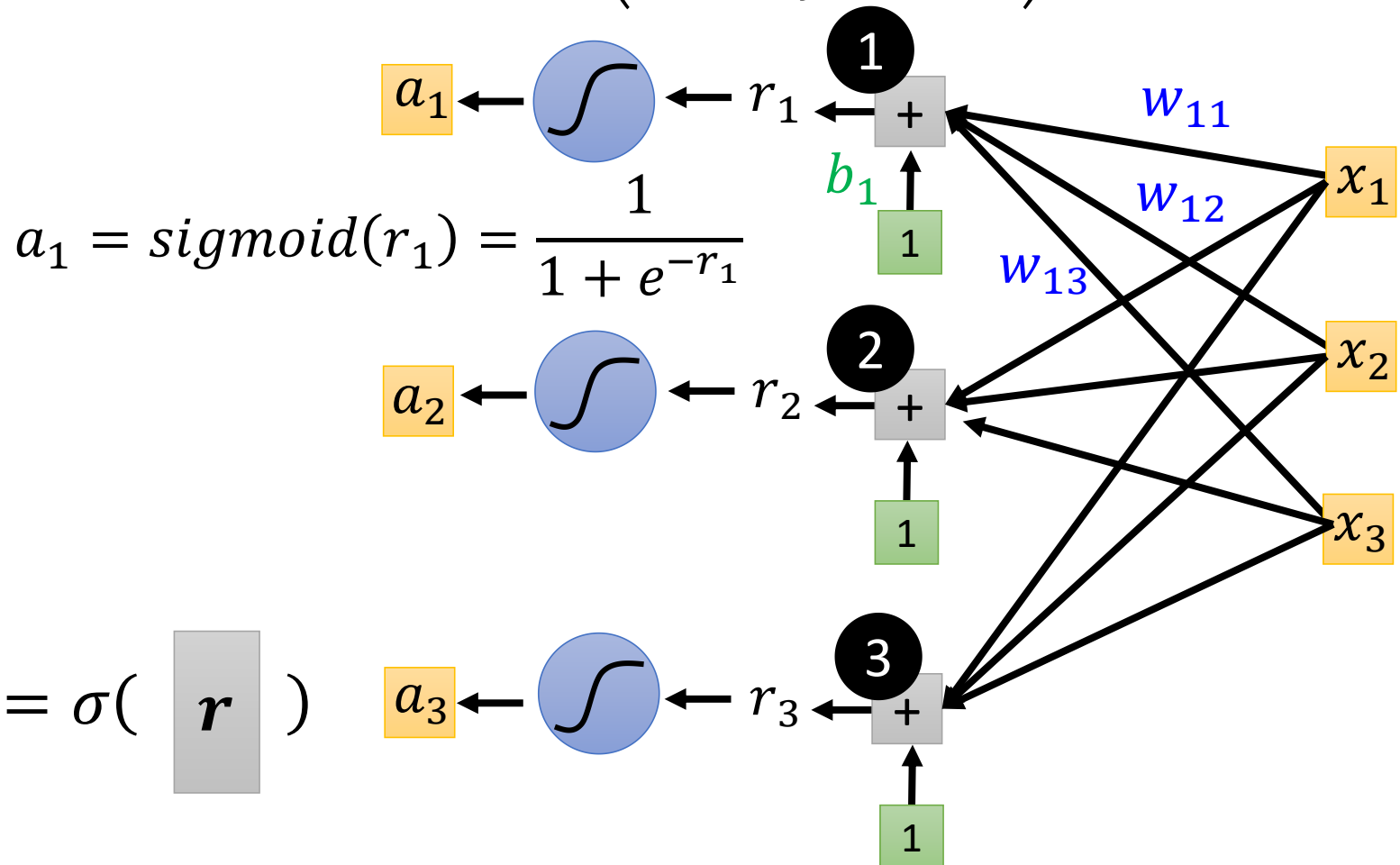
$$\boxed{\textcolor{gray}{r}} = \boxed{\textcolor{green}{b}} + \boxed{\textcolor{blue}{W}} \boxed{\textcolor{orange}{x}}$$

$$y = b + \sum_i \textcolor{red}{c}_i \textit{sigmoid} \left( \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right) \quad \begin{array}{l} i: 1,2,3 \\ j: 1,2,3 \end{array}$$

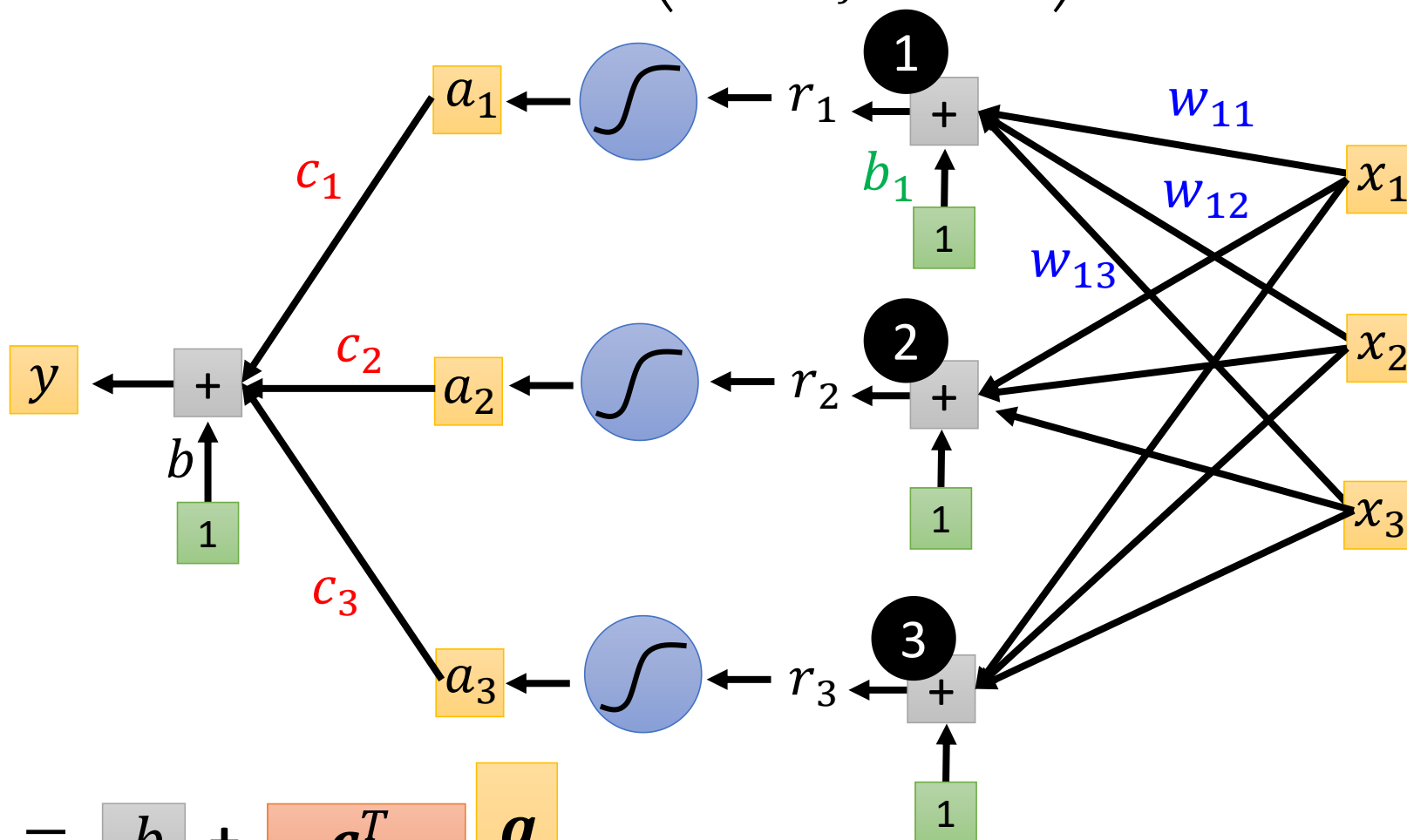
$$\mathbf{r} = \mathbf{b} + \mathbf{W} \mathbf{x}$$



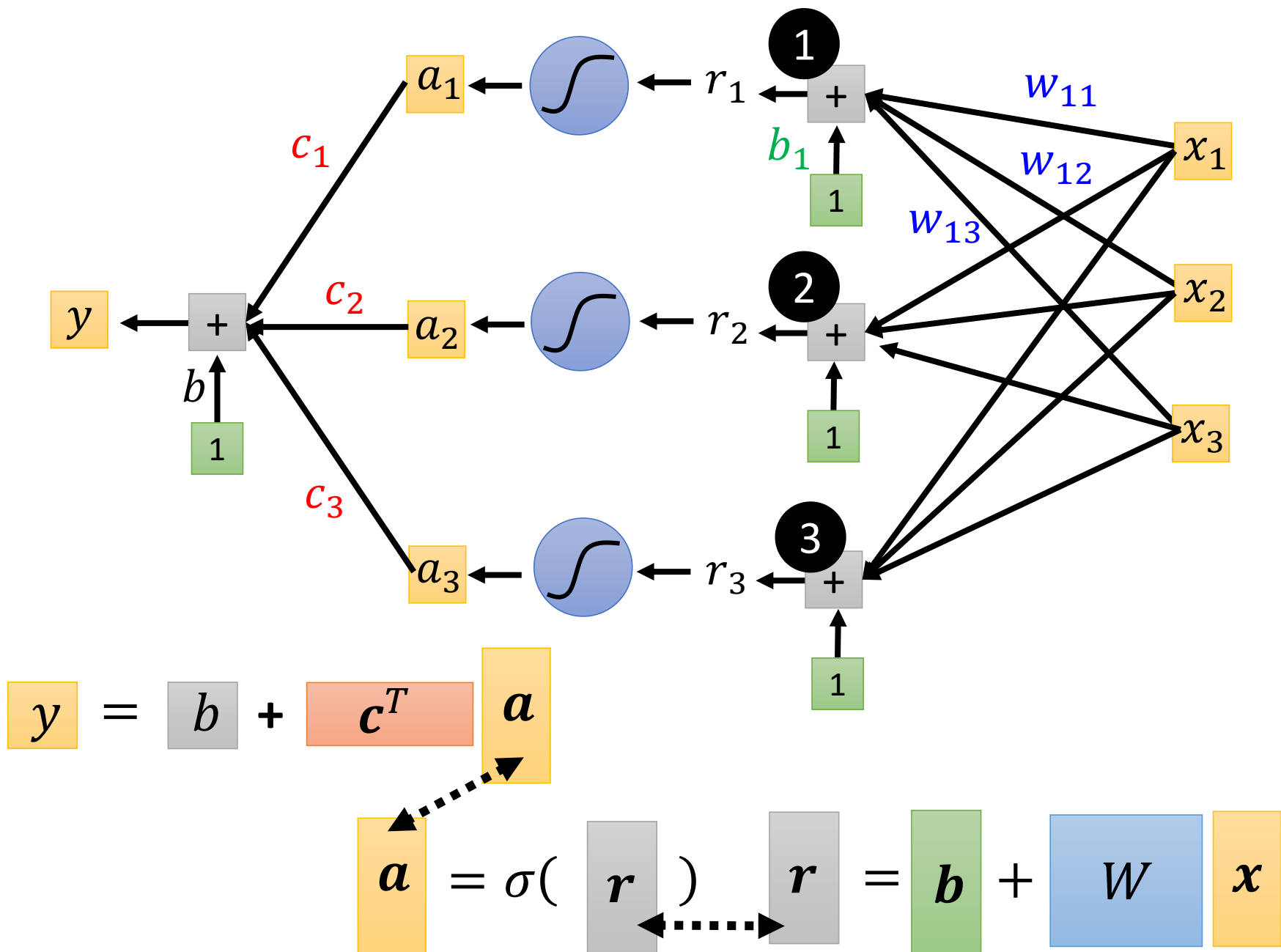
$$y = b + \sum_i \textcolor{red}{c}_i \textcolor{blue}{sigmoid} \left( \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right) \quad \begin{array}{l} i: 1,2,3 \\ j: 1,2,3 \end{array}$$

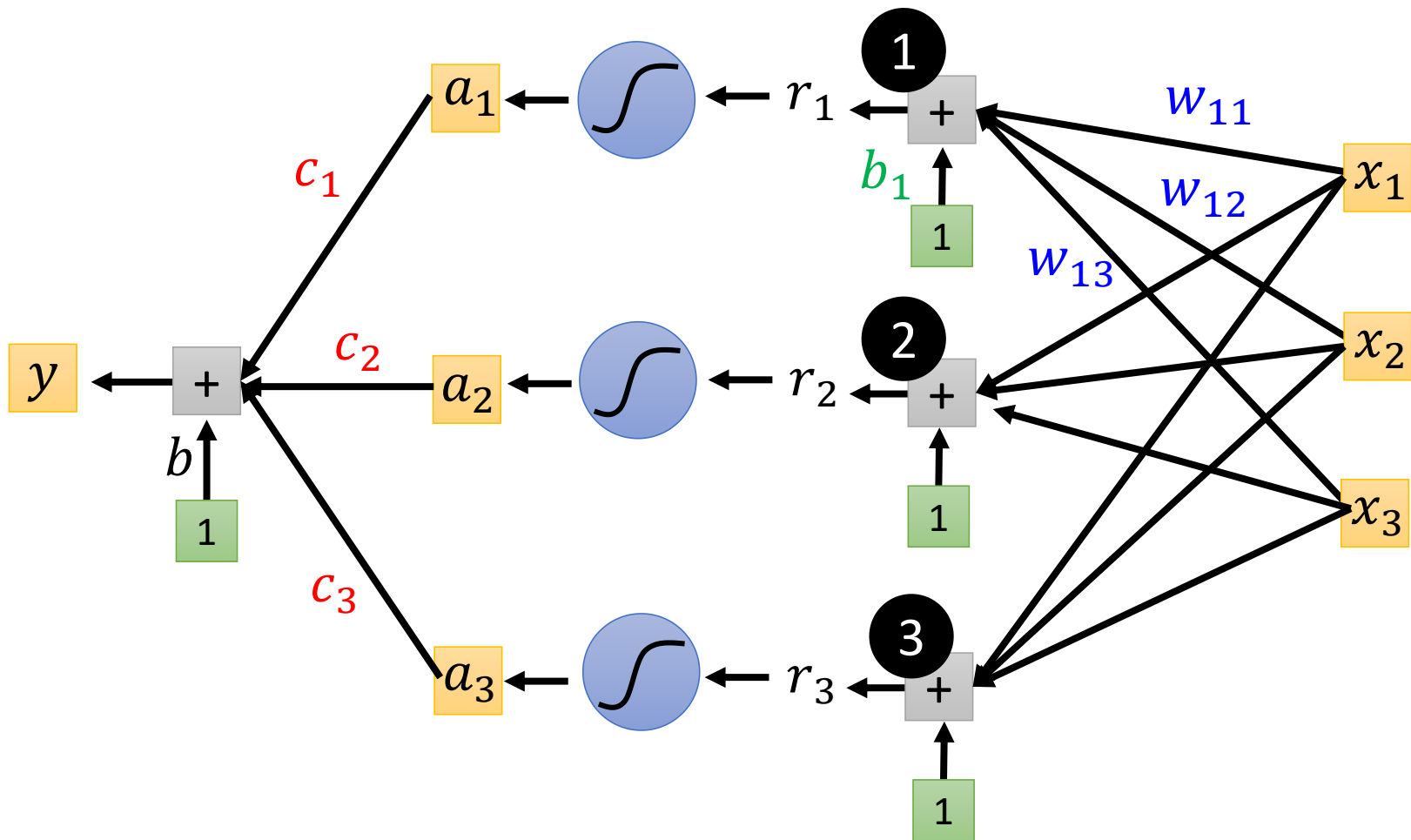


$$y = b + \sum_i \mathbf{c}_i \operatorname{sigmoid} \left( \mathbf{b}_i + \sum_j \mathbf{w}_{ij} x_j \right) \quad \begin{array}{l} i: 1,2,3 \\ j: 1,2,3 \end{array}$$



$$\mathbf{y} = \mathbf{b} + \mathbf{c}^T \mathbf{a}$$





$$y = b + c^T \sigma(b + Wx)$$



# Function with unknown parameters

$$y = b + c^T \sigma( b + W x )$$

$x$  feature

Unknown parameters

$W$

$b$

$c^T$

$b$

Rows of  $W$

$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \end{bmatrix}$

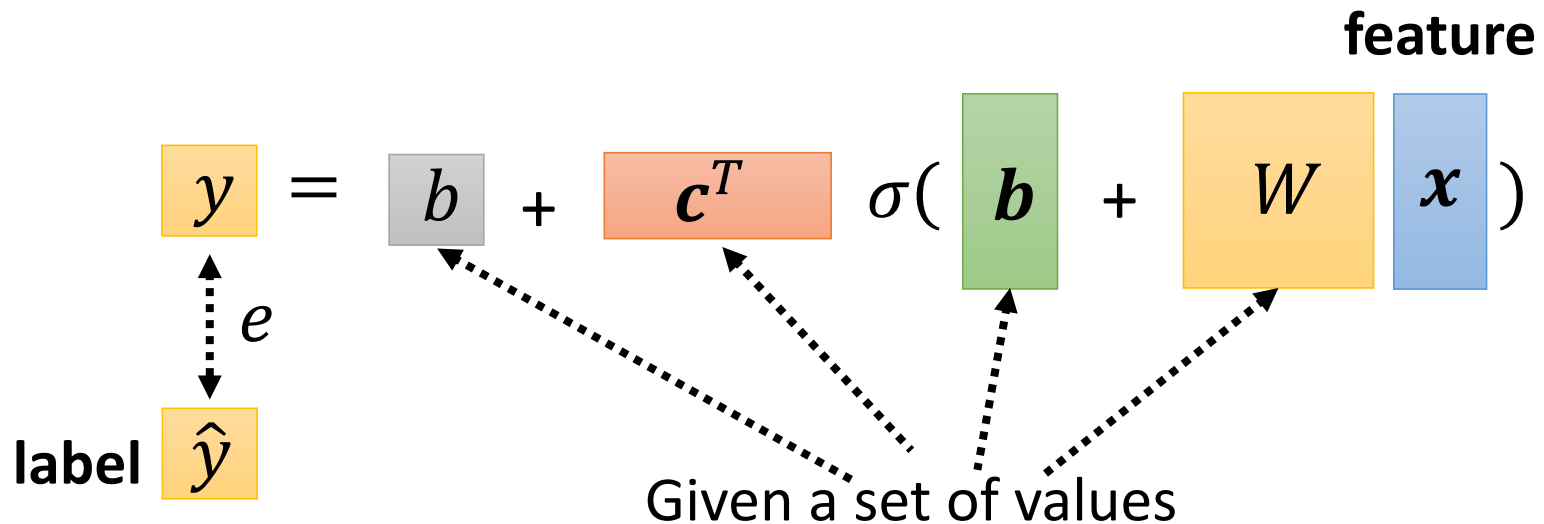
# Back to ML Framework



$$y = b + c^T \sigma( b + W x )$$

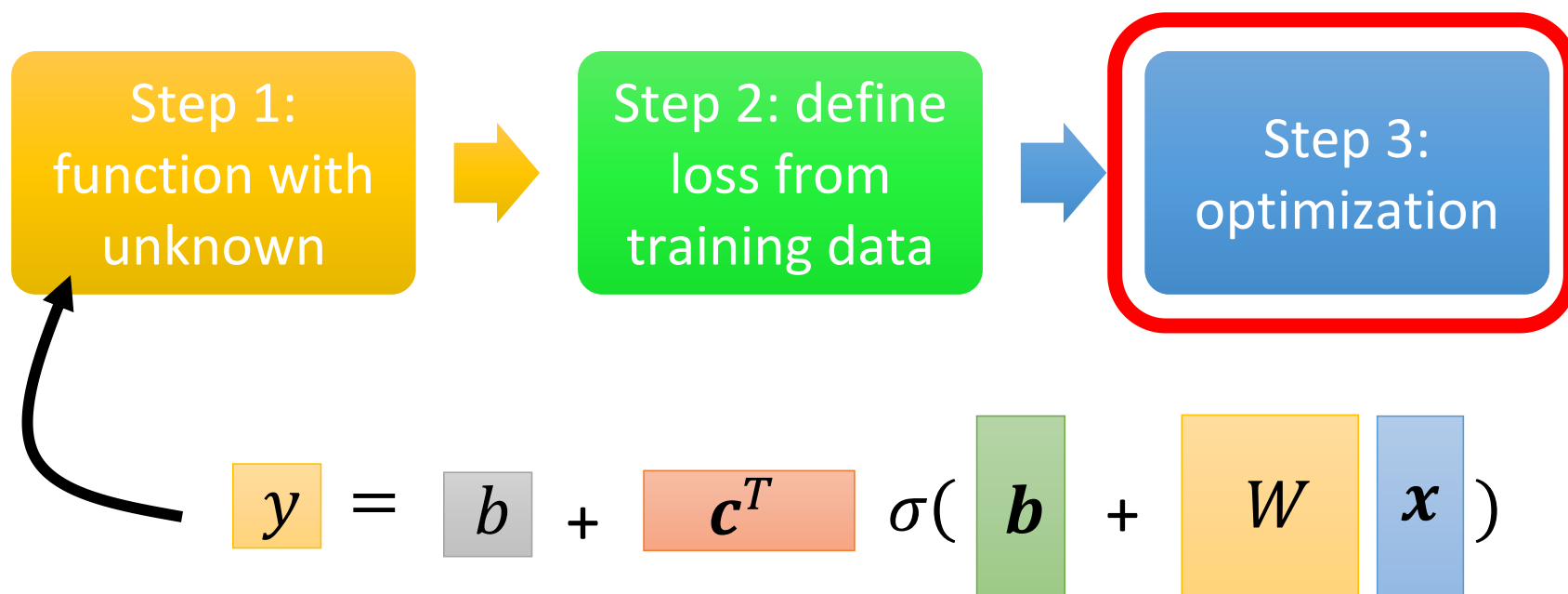
# LOSS

- Loss is a function of parameters  $L(\theta)$
- Loss means how good a set of values is.



Loss: 
$$L = \frac{1}{N} \sum_n e_n$$

# Back to ML Framework



# Optimization of New Model

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} L$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \end{bmatrix}$$

➤ (Randomly) Pick initial values  $\boldsymbol{\theta}^0$

$$\text{gradient } \mathbf{g} = \begin{bmatrix} \left. \frac{\partial L}{\partial \theta_1} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \left. \frac{\partial L}{\partial \theta_2} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix}$$

$$\mathbf{g} = \nabla L(\boldsymbol{\theta}^0)$$

$$\begin{bmatrix} \theta_1^1 \\ \theta_2^1 \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \\ \vdots \end{bmatrix} - \begin{bmatrix} \eta \left. \frac{\partial L}{\partial \theta_1} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \eta \left. \frac{\partial L}{\partial \theta_2} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix}$$

$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \mathbf{g}$$

# Optimization of New Model

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} L$$

➤ (Randomly) Pick initial values  $\boldsymbol{\theta}^0$

➤ Compute gradient  $\boldsymbol{g} = \nabla L(\boldsymbol{\theta}^0)$

$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \boldsymbol{g}$$

➤ Compute gradient  $\boldsymbol{g} = \nabla L(\boldsymbol{\theta}^1)$

$$\boldsymbol{\theta}^2 \leftarrow \boldsymbol{\theta}^1 - \eta \boldsymbol{g}$$

➤ Compute gradient  $\boldsymbol{g} = \nabla L(\boldsymbol{\theta}^2)$

$$\boldsymbol{\theta}^3 \leftarrow \boldsymbol{\theta}^2 - \eta \boldsymbol{g}$$

# Optimization of New Model

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L$$

➤ (Randomly) Pick initial values  $\theta^0$

➤ Compute gradient  $\mathbf{g} = \nabla L^1(\theta^0)$

$$\text{update } \theta^1 \leftarrow \theta^0 - \eta \mathbf{g}$$

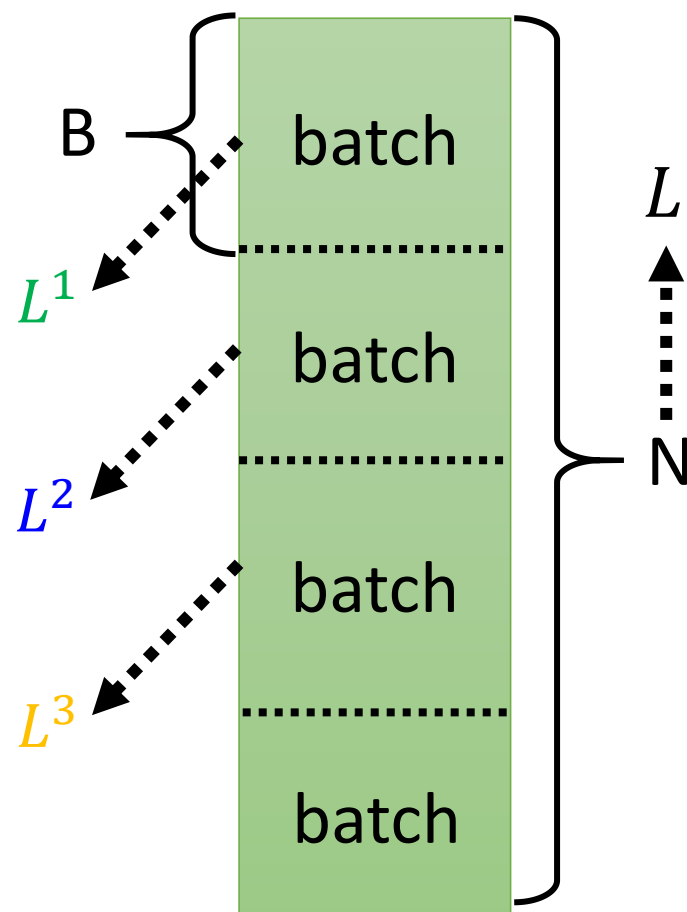
➤ Compute gradient  $\mathbf{g} = \nabla L^2(\theta^1)$

$$\text{update } \theta^2 \leftarrow \theta^1 - \eta \mathbf{g}$$

➤ Compute gradient  $\mathbf{g} = \nabla L^3(\theta^2)$

$$\text{update } \theta^3 \leftarrow \theta^2 - \eta \mathbf{g}$$

1 **epoch** = see all the batches once



# Optimization of New Model

## Example 1

- 10,000 examples ( $N = 10,000$ )
- Batch size is 10 ( $B = 10$ )

How many update in **1 epoch**?

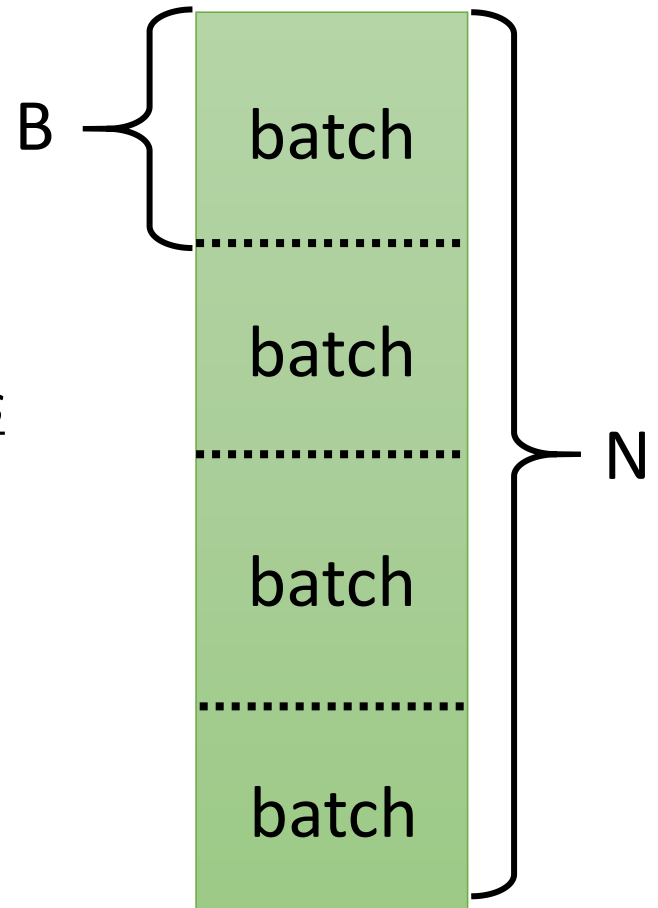
1,000 updates

## Example 2

- 1,000 examples ( $N = 1,000$ )
- Batch size is 100 ( $B = 100$ )

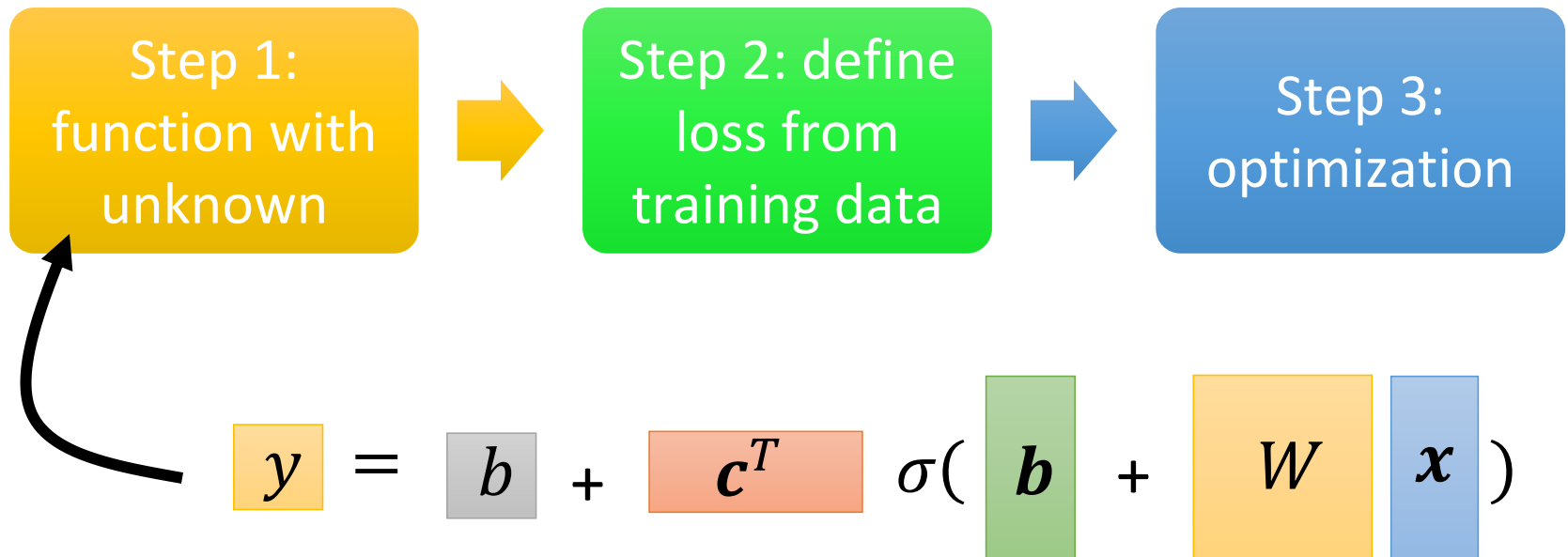
How many update in **1 epoch**?

10 updates





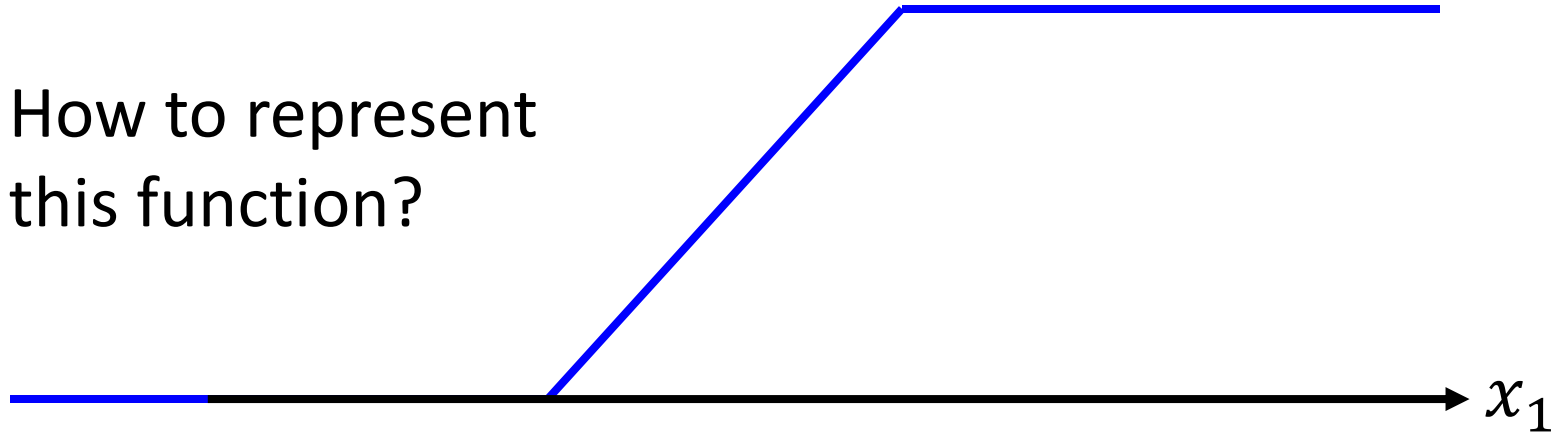
# Back to ML Framework



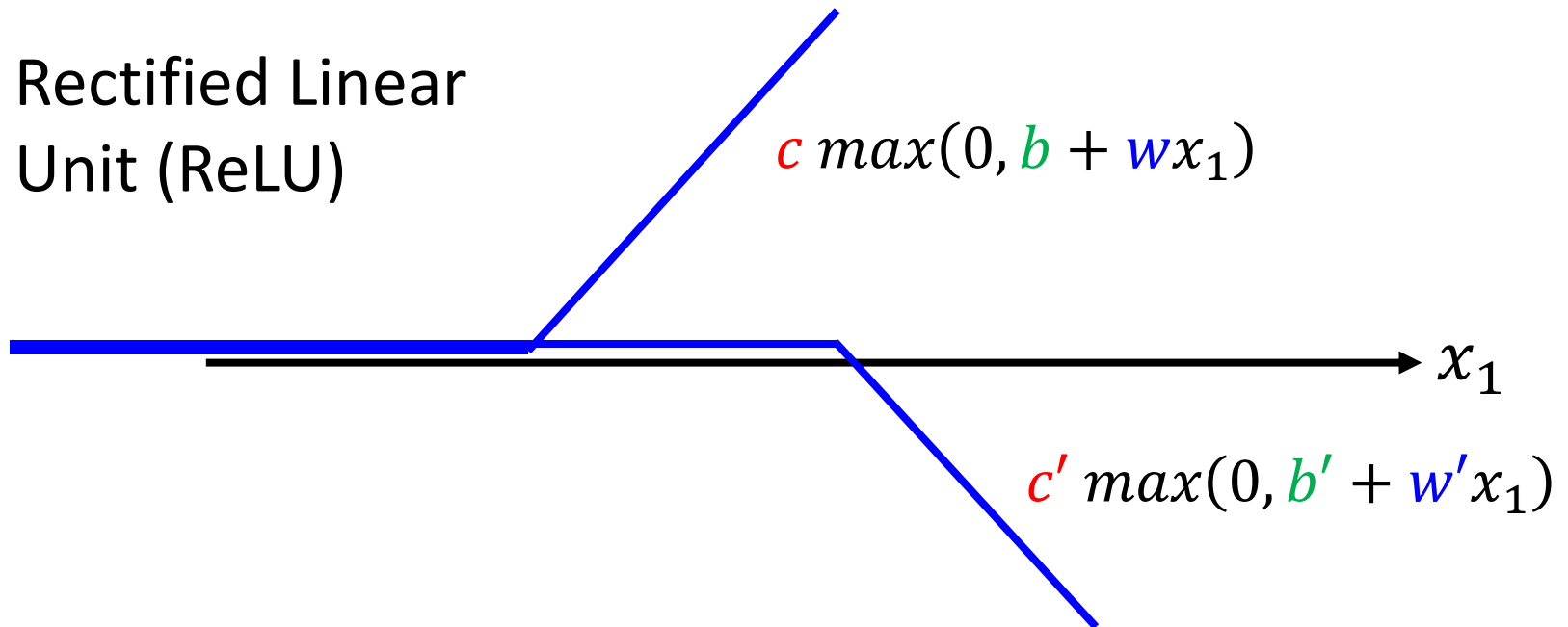
More variety of models ...

# Sigmoid $\rightarrow$ ReLU

How to represent  
this function?



Rectified Linear  
Unit (ReLU)



# Sigmoid → ReLU

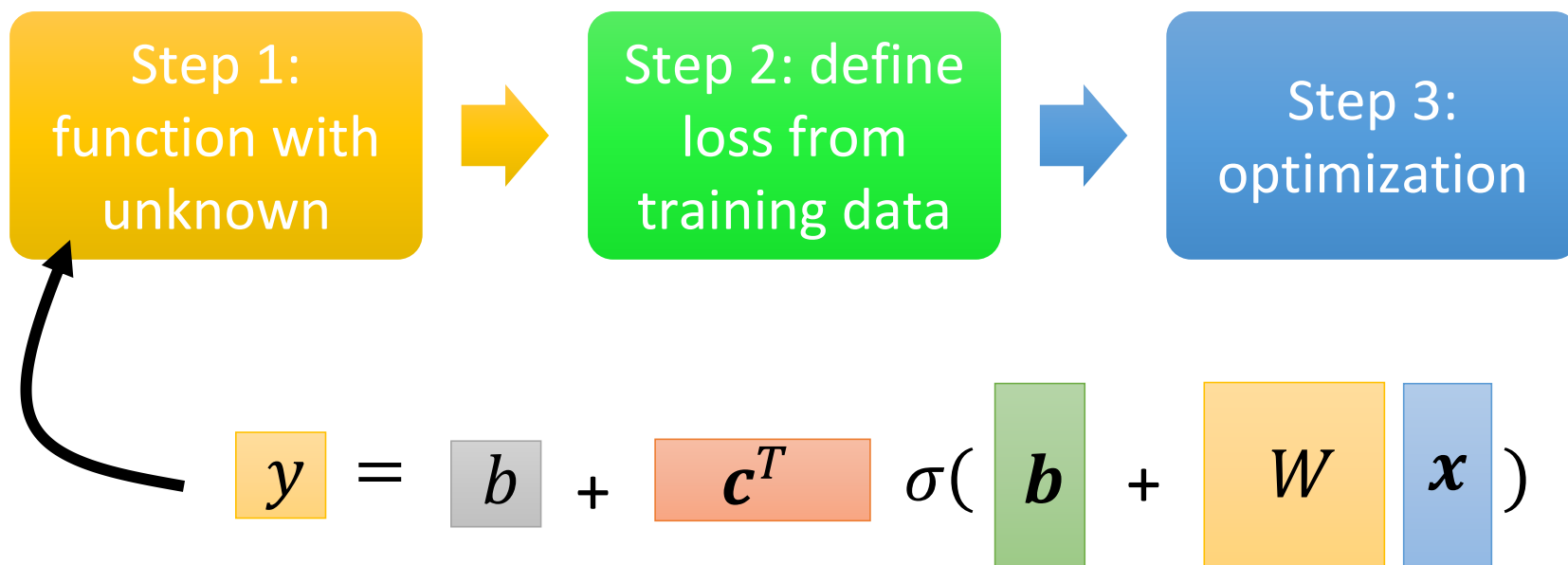
$$y = b + \sum_i c_i \text{sigmoid} \left( b_i + \sum_j w_{ij} x_j \right)$$

**Activation function**

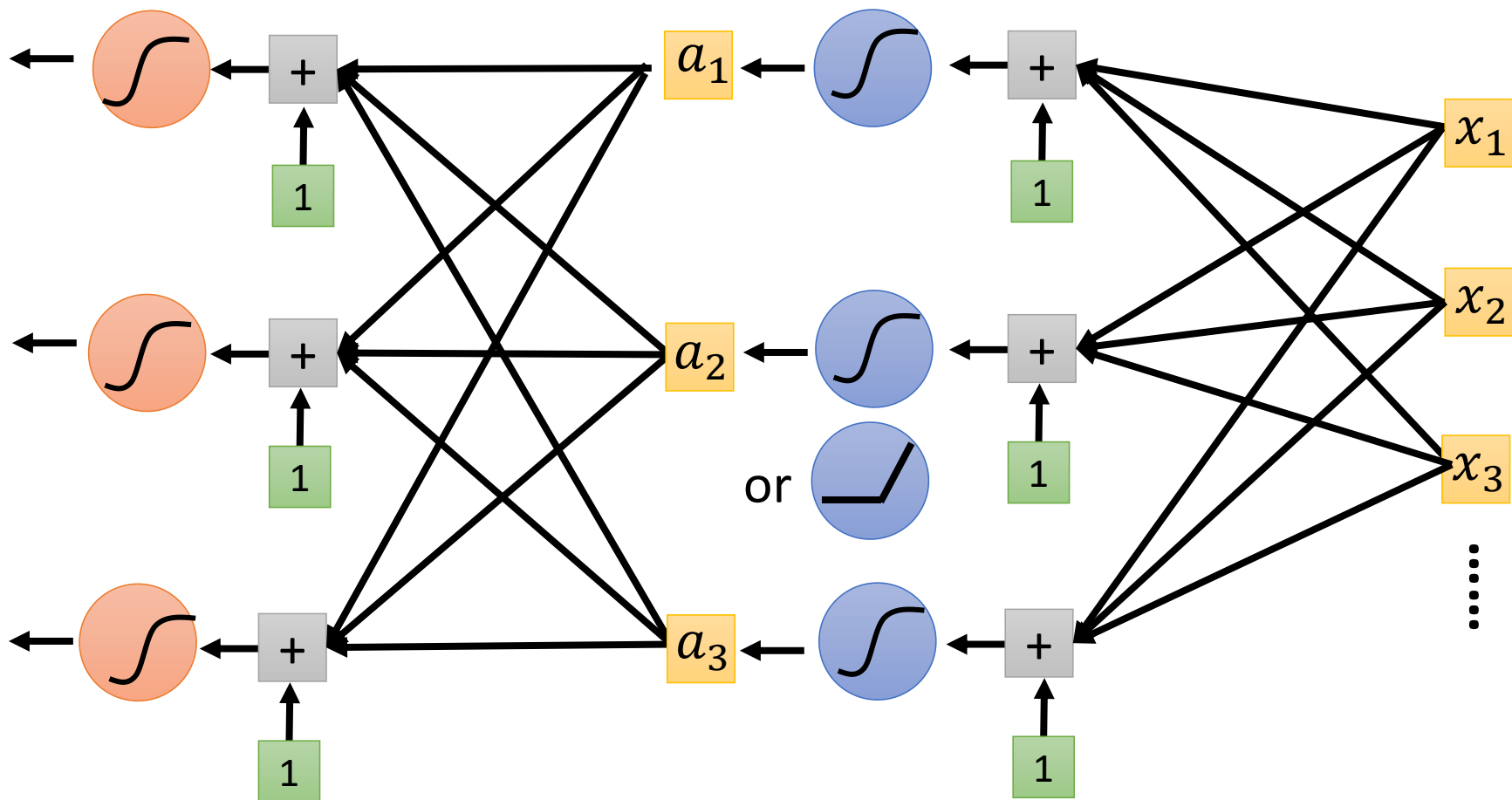
$$y = b + \sum_i c_i \text{max} \left( 0, b_i + \sum_j w_{ij} x_j \right)$$

Which one is better?

# Back to ML Framework



Even more variety of models ...



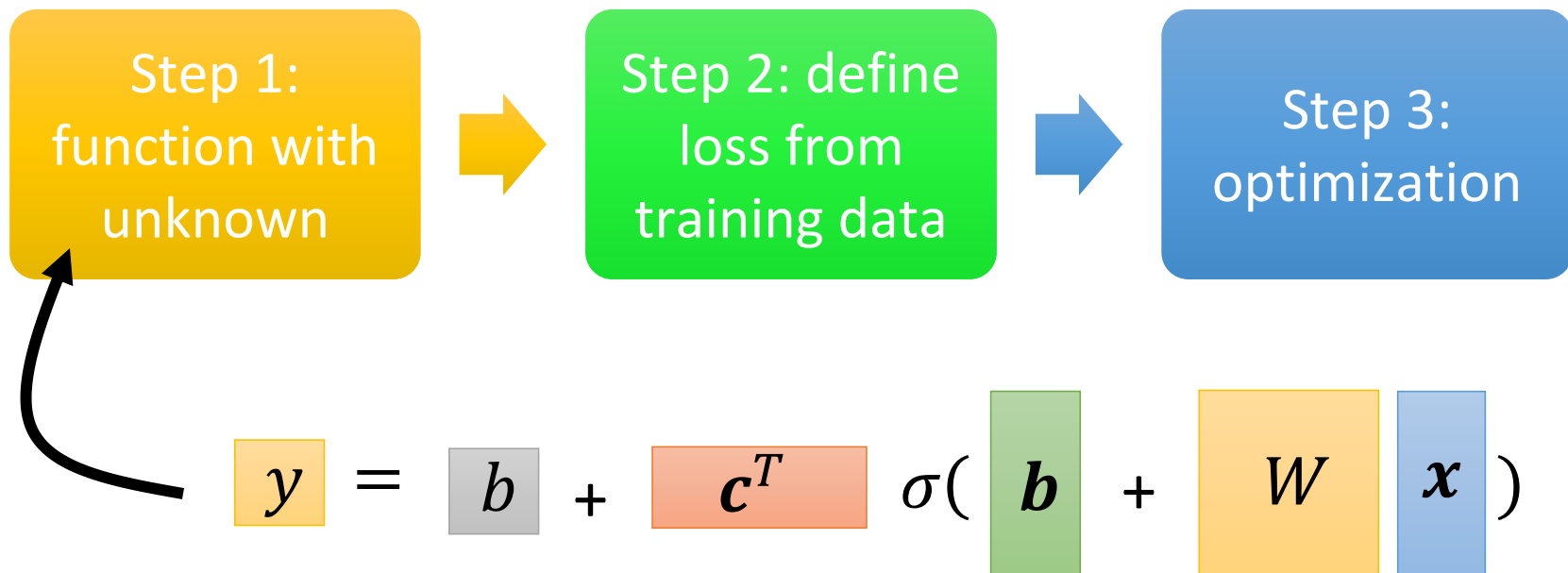
or



$$\begin{aligned}
 \boxed{a'} &= \sigma \left( \boxed{b'} + \boxed{W'} \boxed{a} \right) & \boxed{a} &= \sigma \left( \boxed{b} + \boxed{W} \boxed{x} \right)
 \end{aligned}$$

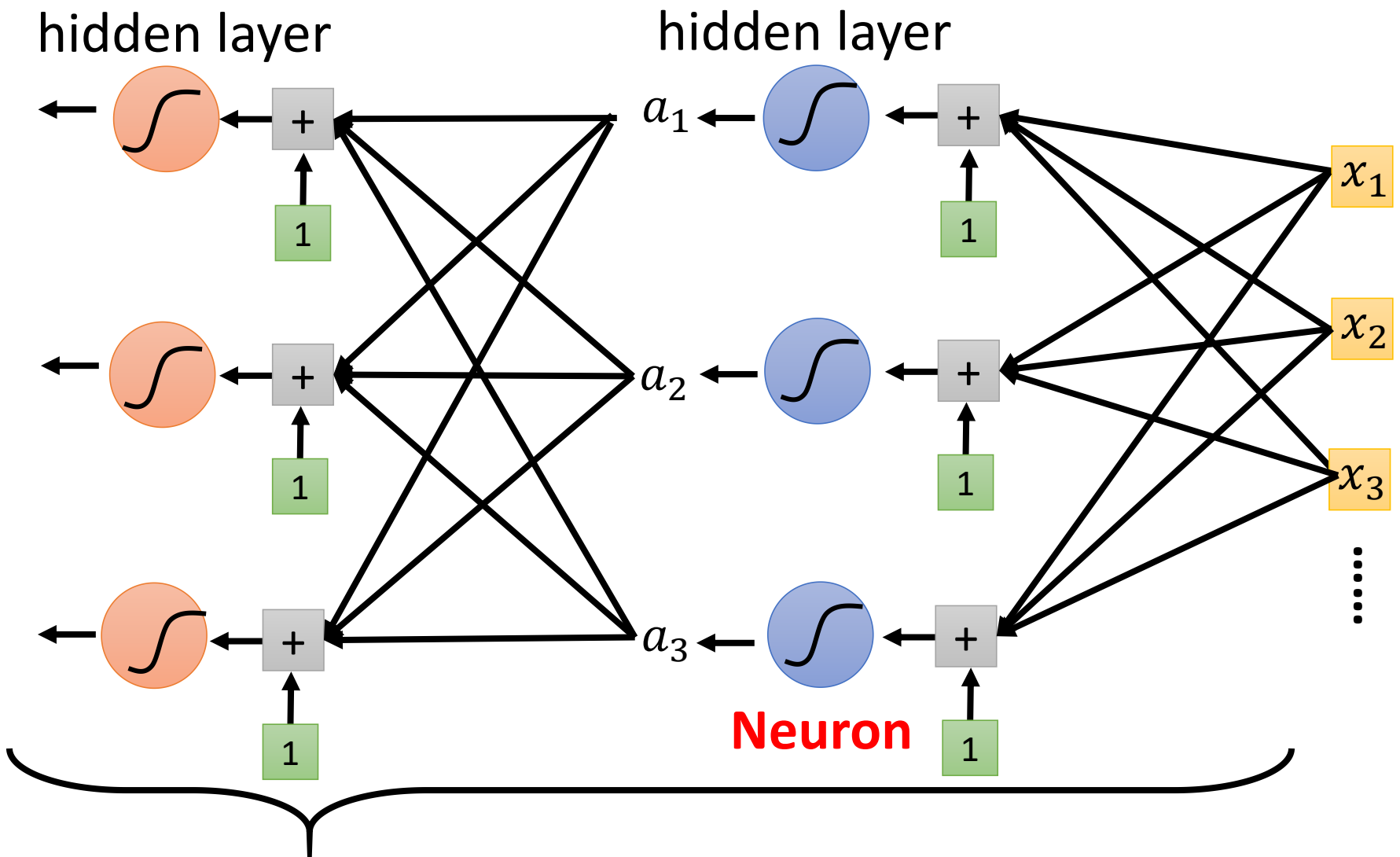
A dashed arrow points from the  $\boxed{a}$  in the first equation to the  $\boxed{a}$  in the second equation, indicating that the output of the first layer is the input to the second layer.

# Back to ML Framework



It is not ***fancy*** enough.

Let's give it a ***fancy*** name!



**Neural Network**

This mimics human brains ... (???)

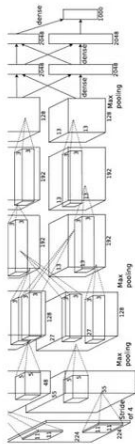
Many layers means **Deep** ➡ Deep Learning

# Deep = Many hidden layers

[http://cs231n.stanford.edu/slides/winter1516\\_lecture8.pdf](http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf)

8 layers

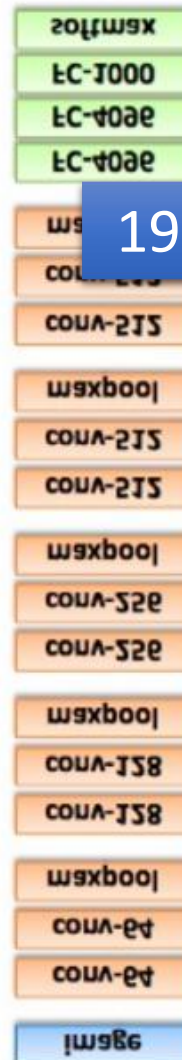
16.4%



AlexNet (2012)

19 layers

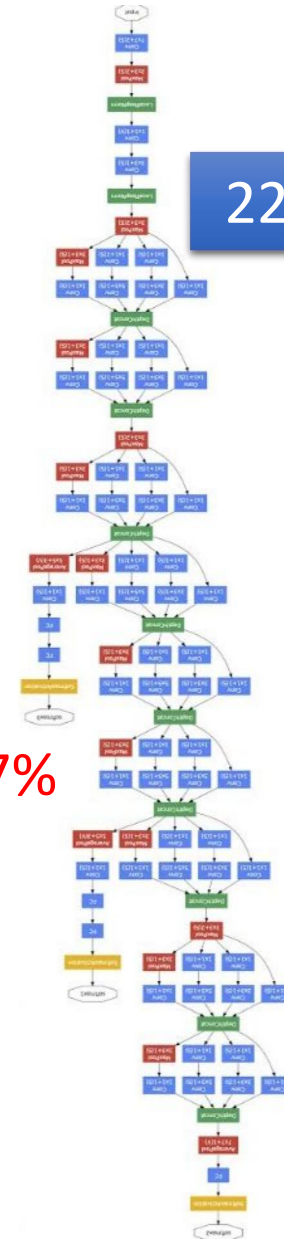
7.3%



VGG (2014)

22 layers

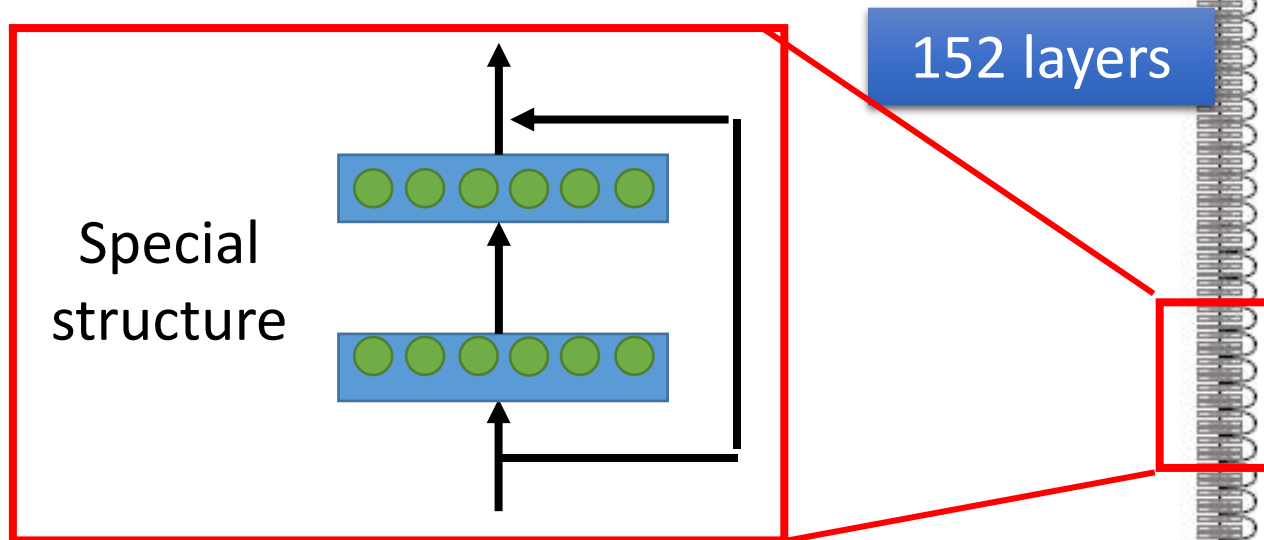
6.7%



GoogleNet (2014)



# Deep = Many hidden layers



Why we want “**Deep**” network,  
not “**Fat**” network?

3.57%

16.4%

AlexNet  
(2012)

7.3%

VGG  
(2014)

6.7%

GoogleNet  
(2014)

Residual Net  
(2015)