

# Introduction of Machine Learning

---

# Machine Learning $\approx$ Looking for Function

- Speech Recognition

$$f\left(\text{[Waveform]}\right) = \text{"How are you"}$$

- Image Recognition

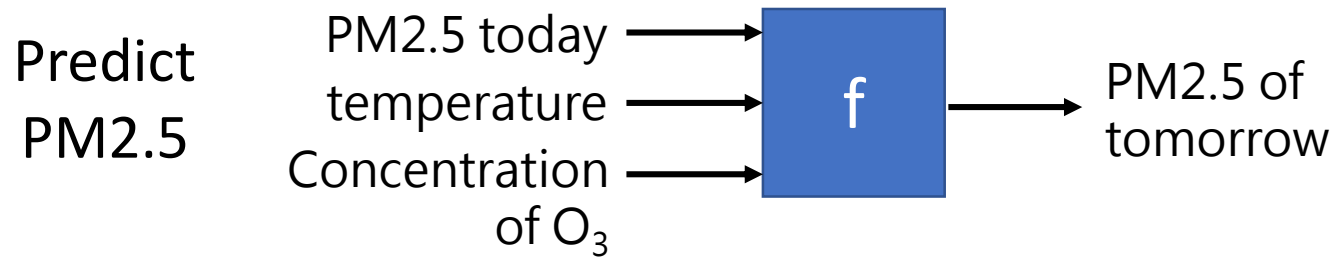
$$f\left(\text{[Cat Image]}\right) = \text{"Cat"}$$

- Playing Go

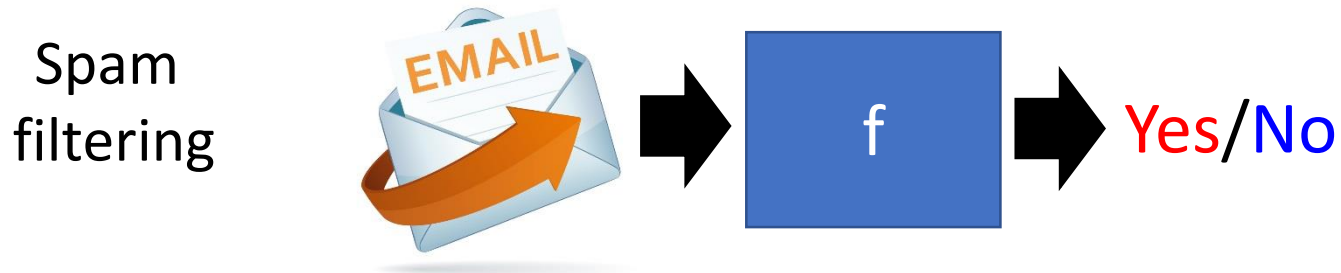
$$f\left(\text{[Go Board Image]}\right) = \text{"5-5"}_{\text{(next move)}}$$

# Different types of Functions

**Regression**: The function outputs a scalar.



**Classification**: Given options (**classes**), the function outputs the correct one.

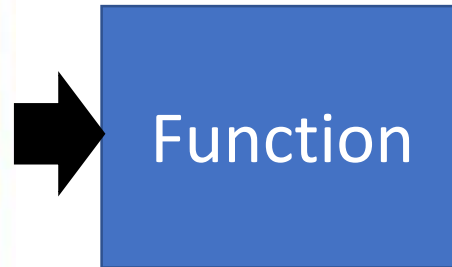


# Different types of Functions

**Classification:** Given options (**classes**), the function outputs the correct one.



**Playing GO**



a position on  
the board

**Next move**

Each position  
is a class  
(19 x 19 classes)

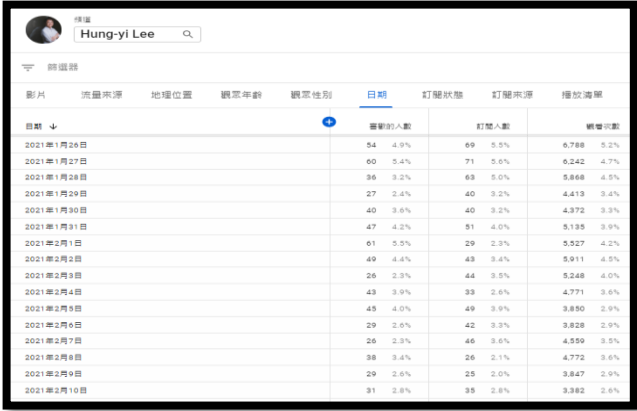


How to find a function?  
A Case Study

# The function we want to find ...

$$y = f(\quad)$$

no. of views  
on 2/26



The screenshot shows a YouTube channel page for 'Hung-yi Lee'. Below the channel header, there is a navigation bar with tabs: '影片' (Videos), '流量來源' (Traffic sources), '地理位置' (Geography), '觀眾年齡' (Audience age), '觀眾性別' (Audience gender), '日期' (Date), '訂閱狀態' (Subscription status), '訂閱來源' (Subscription source), and '播放清單' (Playlists). The '日期' (Date) tab is selected. Below the navigation bar is a table with columns: '日期' (Date), '觀看人數' (Views), '訂閱人數' (Subscribers), and '觀看次數' (Views). The table contains data for various dates from 2021年1月26日 to 2021年3月10日.

日期	觀看人數	訂閱人數	觀看次數
2021年1月26日	54	69	6,788
2021年1月27日	60	71	6,242
2021年1月28日	36	63	5,868
2021年1月29日	27	40	4,413
2021年1月30日	40	40	4,372
2021年1月31日	47	51	5,135
2021年2月1日	61	29	5,527
2021年2月2日	49	43	5,911
2021年2月3日	26	44	5,248
2021年2月4日	43	33	4,771
2021年2月5日	45	49	3,850
2021年2月6日	29	42	3,828
2021年2月7日	20	46	4,559
2021年2月8日	38	26	4,772
2021年2月9日	29	25	3,847
2021年3月10日	31	35	3,382



# 1. Function with Unknown Parameters

$$y = f($$



**Model**  $y = b + wx_1$  based on domain knowledge

**feature**

$y$ : no. of views on 2/26,  $x_1$ : no. of views on 2/25

$w$  and  $b$  are unknown parameters (learned from data)

**weight**    **bias**

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

Hung-yi Lee

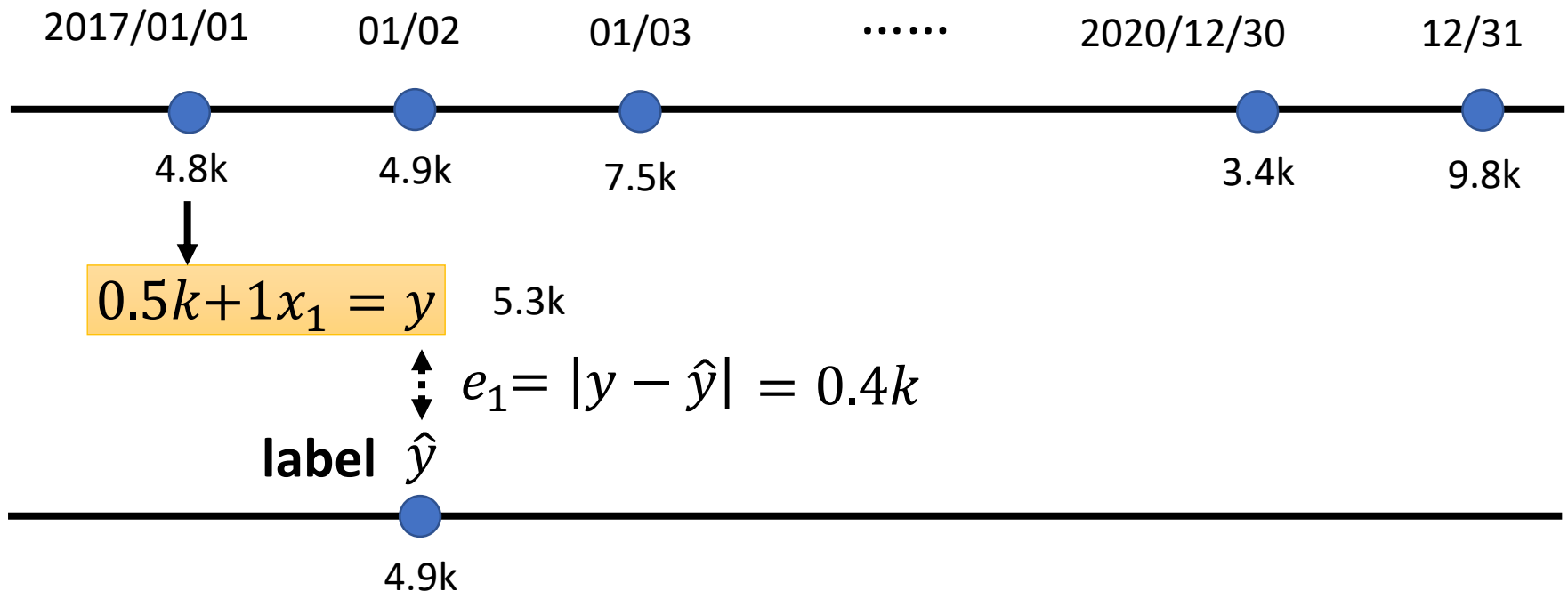
Hung-yi Lee

## 2. Define **Loss** from Training Data

- Loss is a function of parameters  $L(b, w)$
- Loss: how good a set of values is.

$L(0.5k, 1)$   $y = b + wx_1 \longrightarrow y = 0.5k + 1x_1$  How good it is?

Data from 2017/01/01 – 2020/12/31



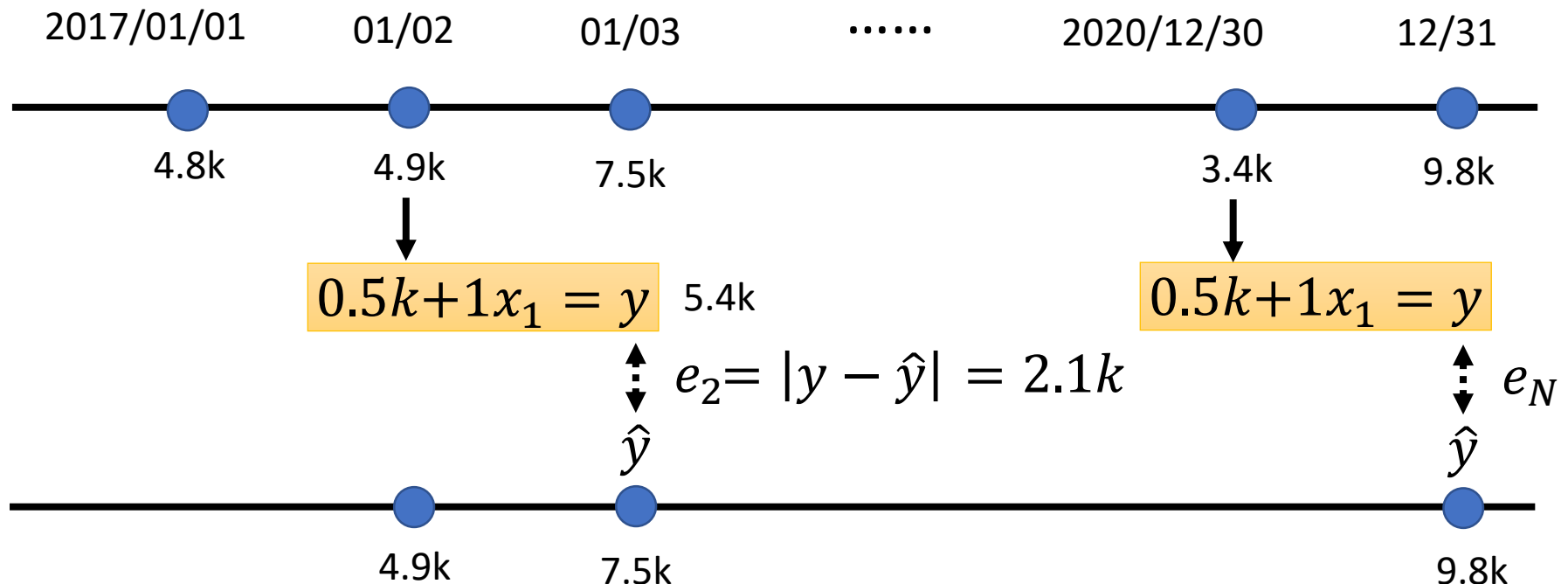


## 2. Define **Loss** from Training Data

- Loss is a function of parameters  $L(b, w)$
- Loss: how good a set of values is.

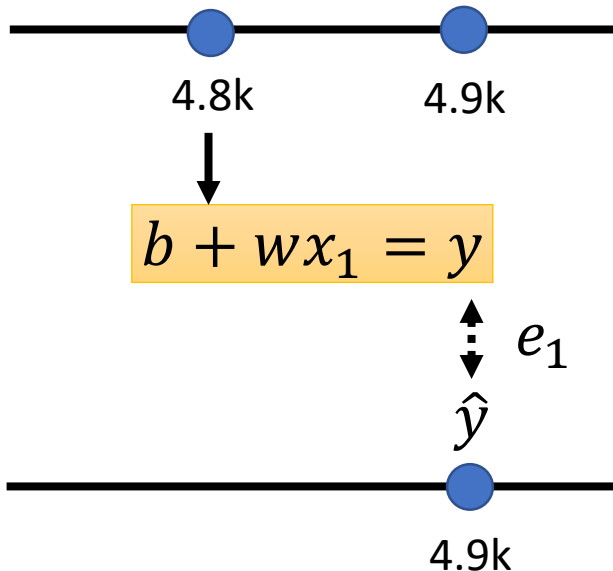
$L(0.5k, 1)$   $y = b + wx_1 \longrightarrow y = 0.5k + 1x_1$  How good it is?

Data from 2017/01/01 – 2020/12/31



## 2. Define Loss from Training Data

- Loss is a function of parameters  $L(b, w)$
- Loss: how good a set of values is.



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

$e = |y - \hat{y}|$        $L$  is mean absolute error (MAE)

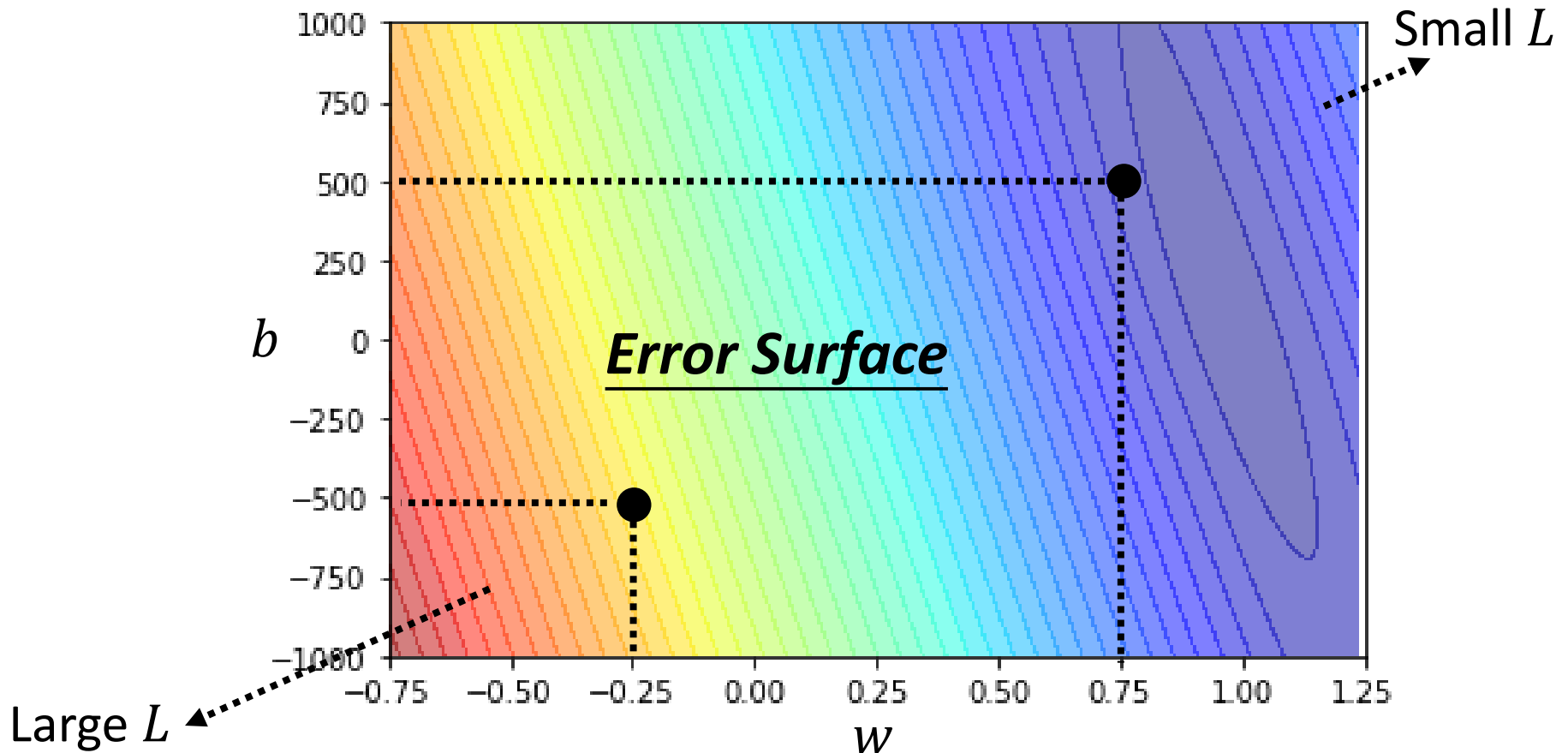
$e = (y - \hat{y})^2$        $L$  is mean square error (MSE)

If  $y$  and  $\hat{y}$  are both probability distributions ➡ Cross-entropy

## 2. Define Loss from Training Data

- Loss is a function of parameters  $L(b, w)$
- Loss: how good a set of values is.

**Model**  $y = b + wx_1$

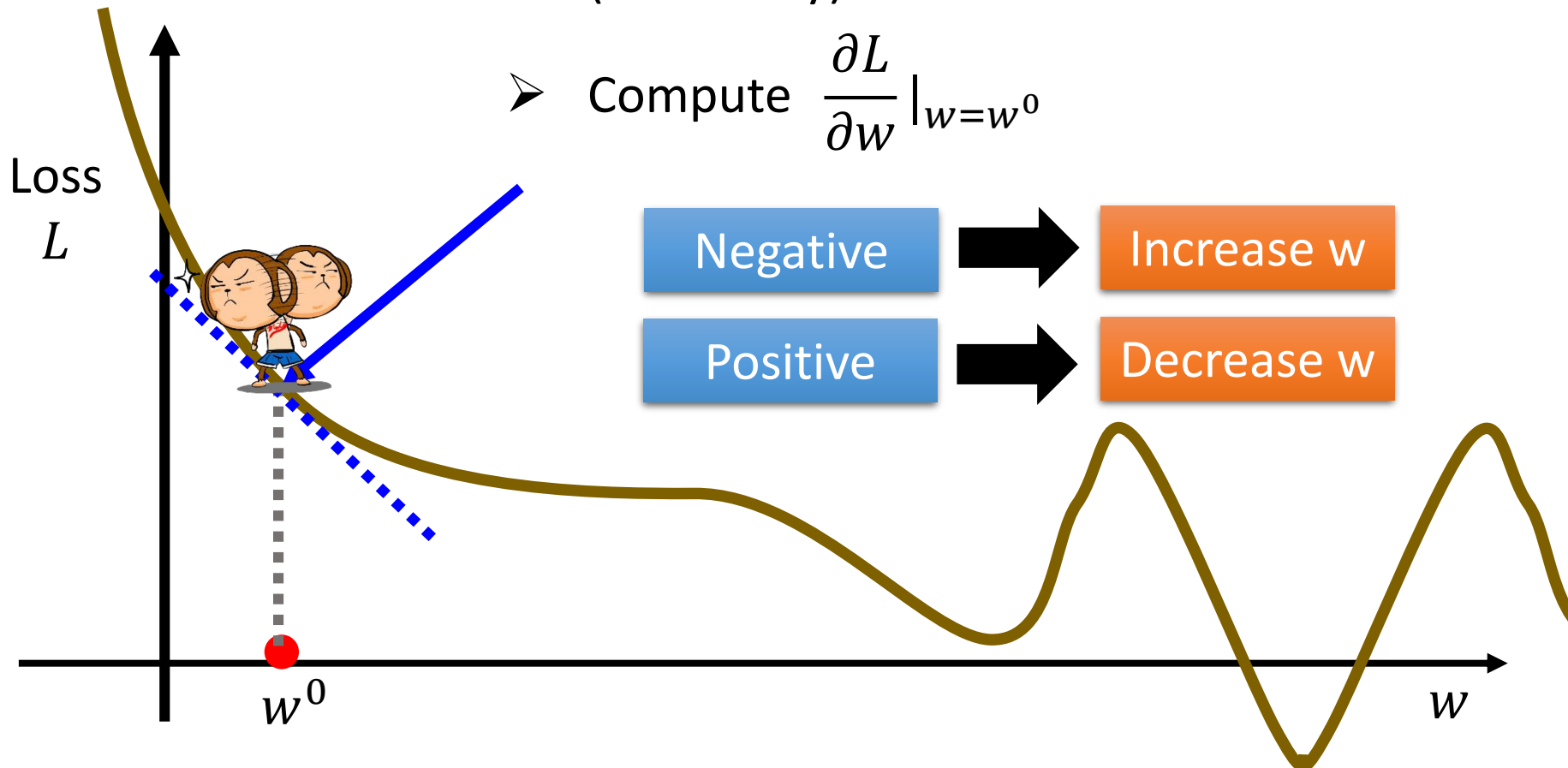


# 3. Optimization

$$w^* = \arg \min_w L$$

## Gradient Descent

- (Randomly) Pick an initial value  $w^0$
- Compute  $\frac{\partial L}{\partial w} \big|_{w=w^0}$



# 3. Optimization

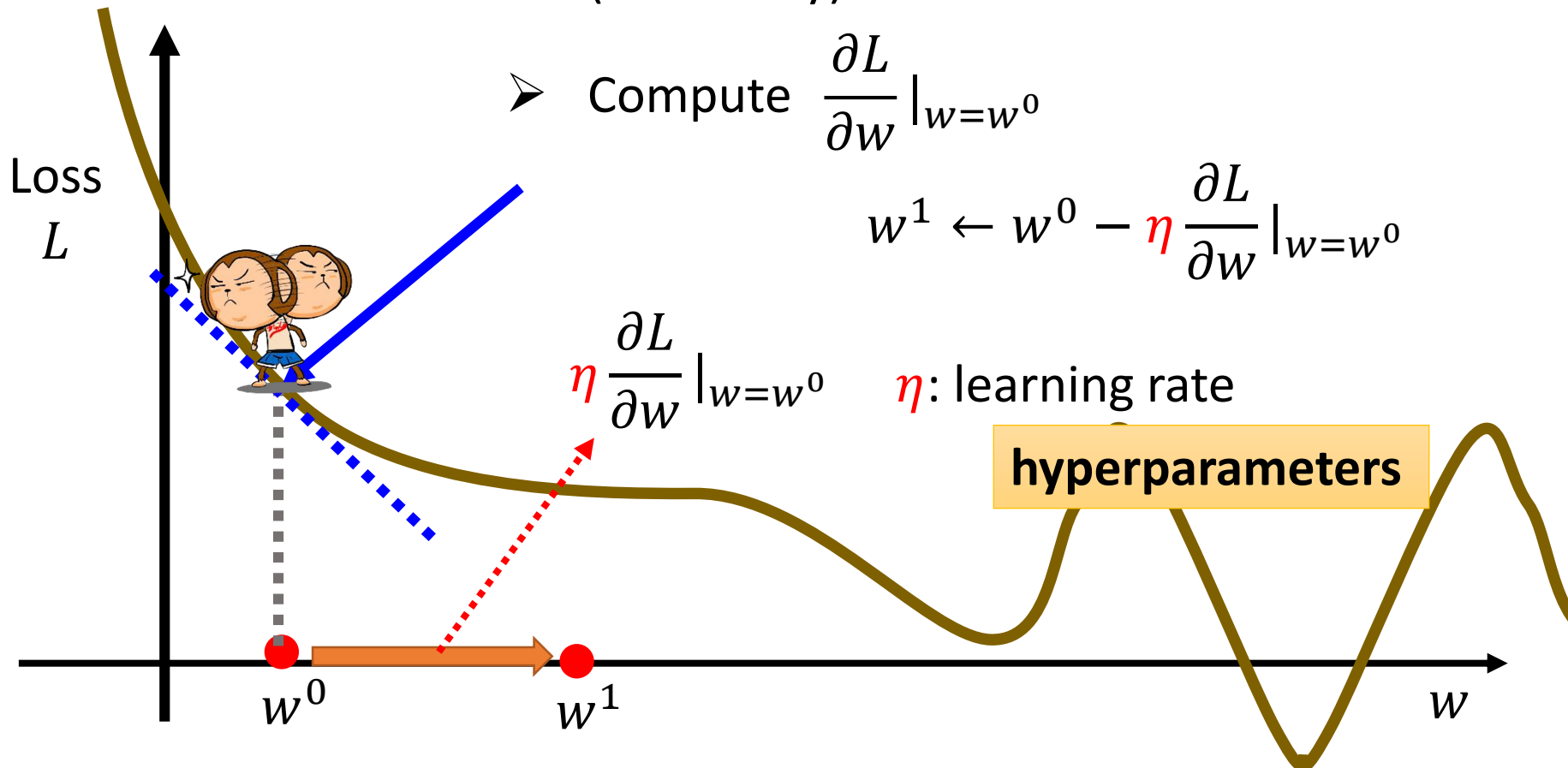
$$w^* = \arg \min_w L$$

## Gradient Descent

➤ (Randomly) Pick an initial value  $w^0$

➤ Compute  $\frac{\partial L}{\partial w} \big|_{w=w^0}$

$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w} \big|_{w=w^0}$$



# 3. Optimization

$$w^* = \arg \min_w L$$

## Gradient Descent

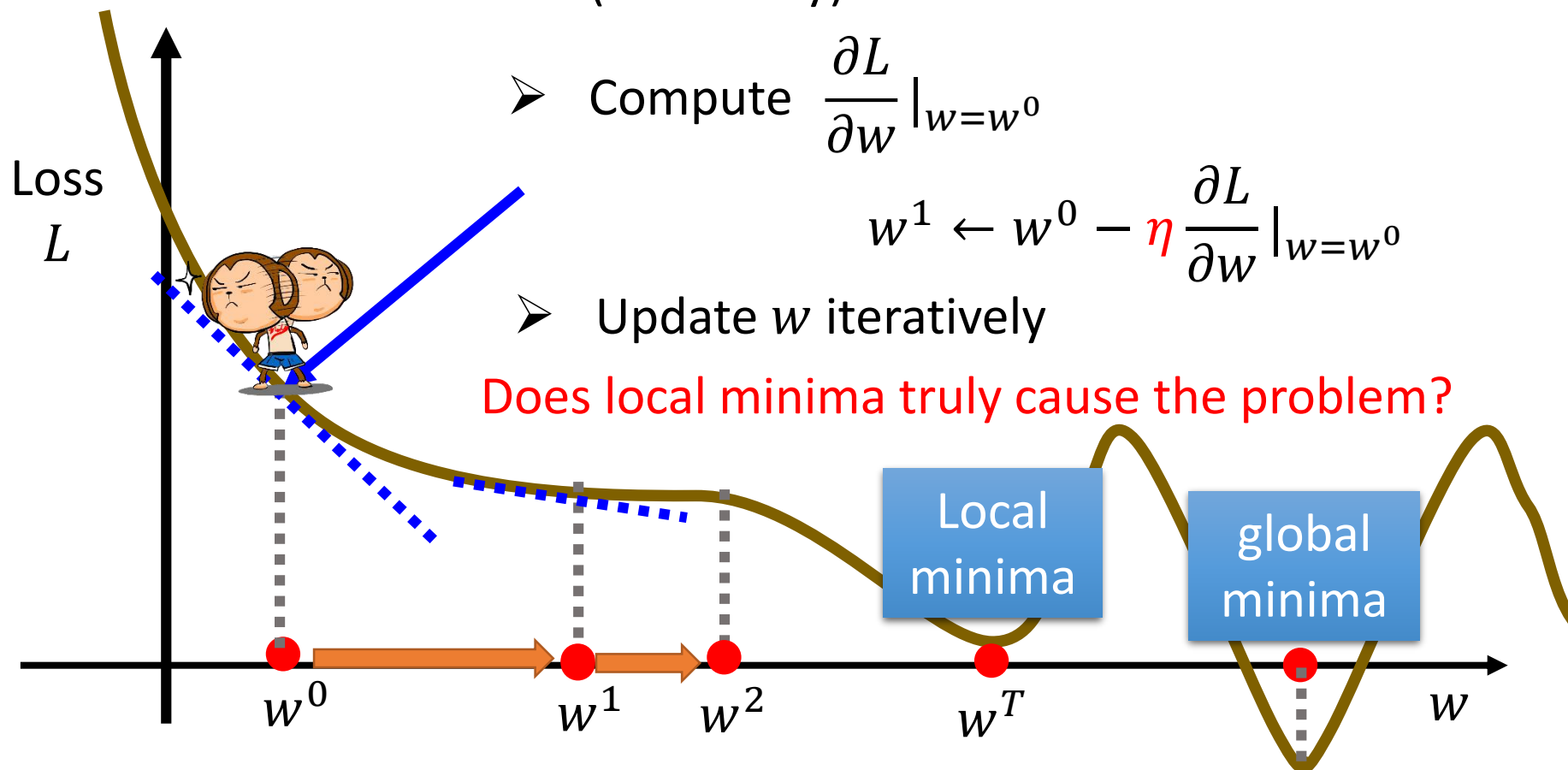
➤ (Randomly) Pick an initial value  $w^0$

➤ Compute  $\frac{\partial L}{\partial w} \big|_{w=w^0}$

$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w} \big|_{w=w^0}$$

➤ Update  $w$  iteratively


Does local minima truly cause the problem?



### 3. Optimization

$$w^*, b^* = \arg \min_{w, b} L$$

- (Randomly) Pick initial values  $w^0, b^0$
- Compute

$$\begin{aligned} \frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0} \\ \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0} \end{aligned}$$


$$w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0}$$

$$b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$$

Can be done in one line in most deep learning frameworks

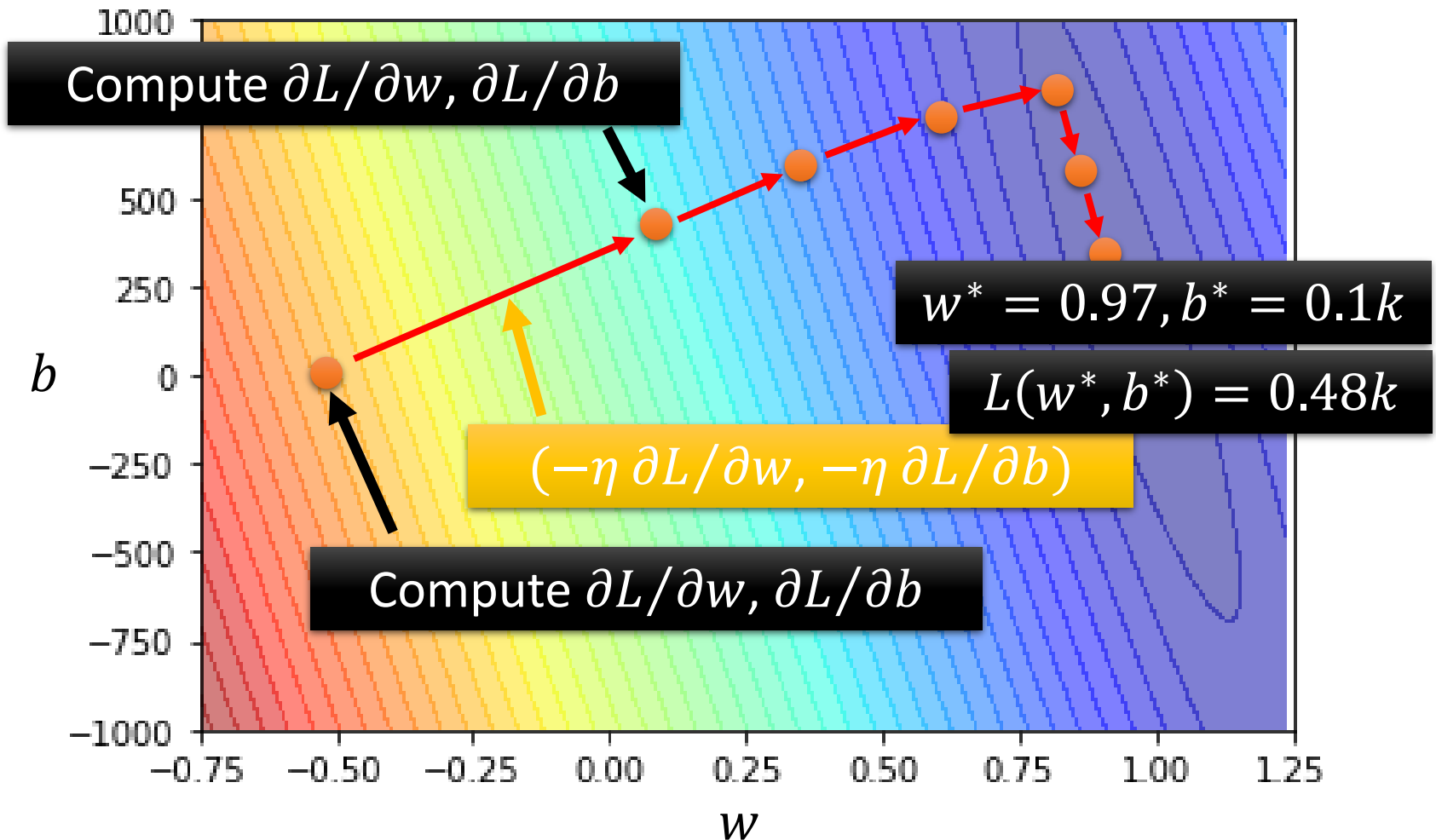
- Update  $w$  and  $b$  iteratively



Model  $y = b + wx_1$

### 3. Optimization

$$w^*, b^* = \arg \min_{w, b} L$$



# Machine Learning is so simple .....

$$w^* = 0.97, b^* = 0.1k$$

$$L(w^*, b^*) = 0.48k$$

$$y = b + wx_1$$

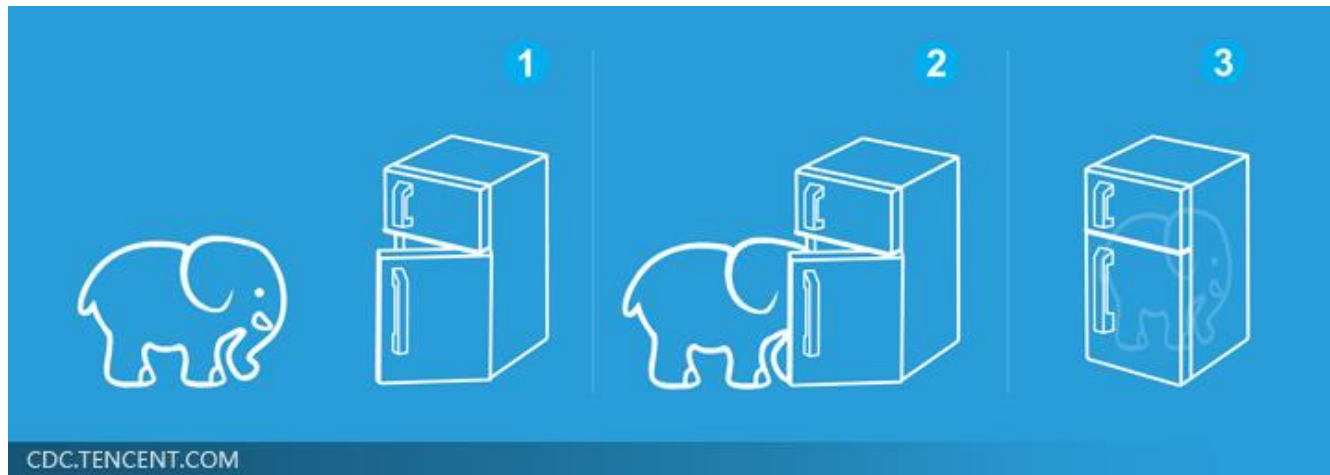
Step 1:  
function with  
unknown



Step 2: define  
loss from  
training data



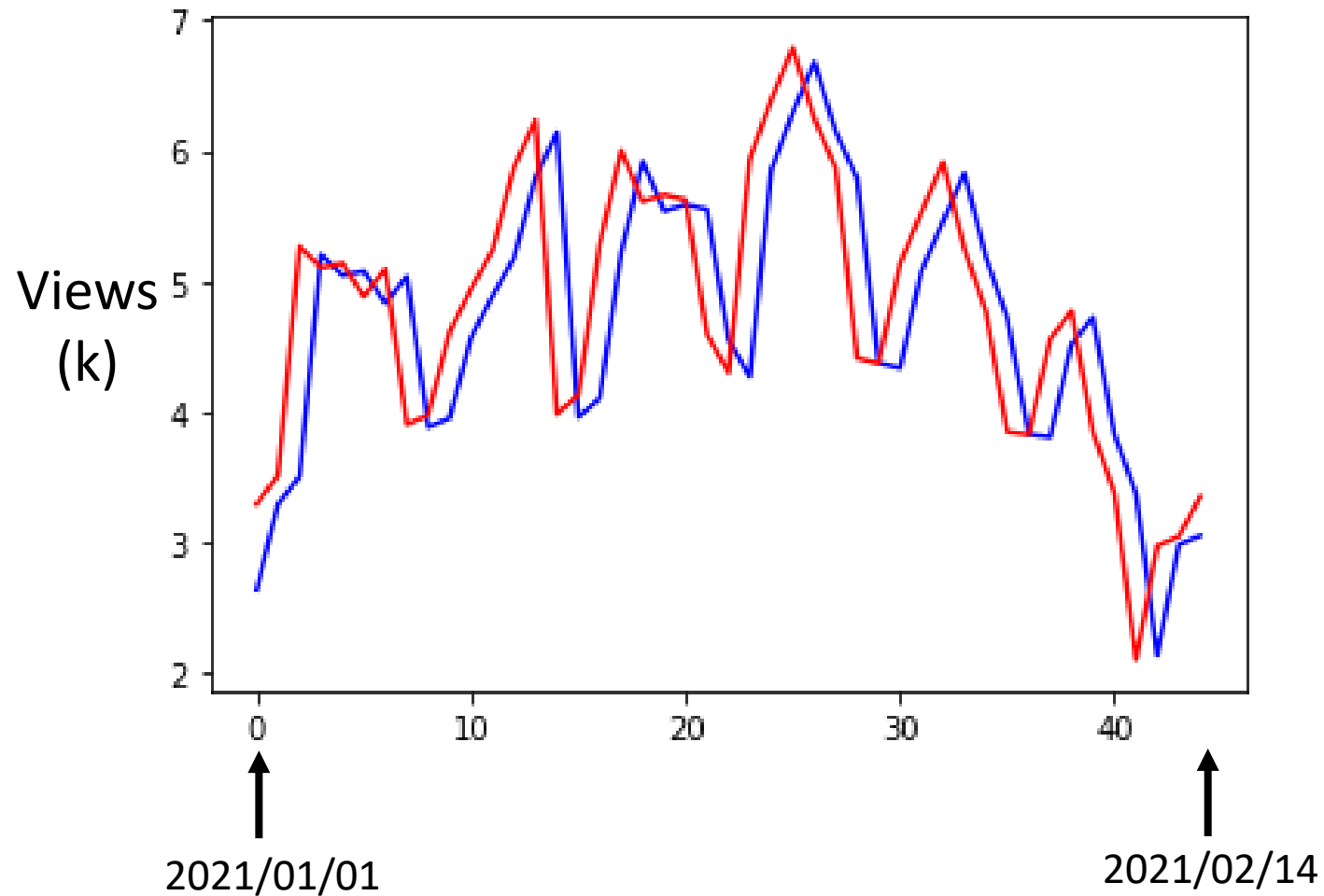
Step 3:  
optimization



$$y = 0.1k + 0.97x_1$$

Red: real no. of views

blue: estimated no. of views



$$y = b + wx_1$$

2017 - 2020

$$L = 0.48k$$

2021

$$L' = 0.58k$$

$$y = b + \sum_{j=1}^7 w_j x_j$$

2017 - 2020

$$L = 0.38k$$

2021

$$L' = 0.49k$$

$b$	$w_1^*$	$w_2^*$	$w_3^*$	$w_4^*$	$w_5^*$	$w_6^*$	$w_7^*$
0.05k	0.79	-0.31	0.12	-0.01	-0.10	0.30	0.18

$$y = b + \sum_{j=1}^{28} w_j x_j$$

2017 - 2020

$$L = 0.33k$$

2021

$$L' = 0.46k$$

$$y = b + \sum_{j=1}^{56} w_j x_j$$

2017 - 2020

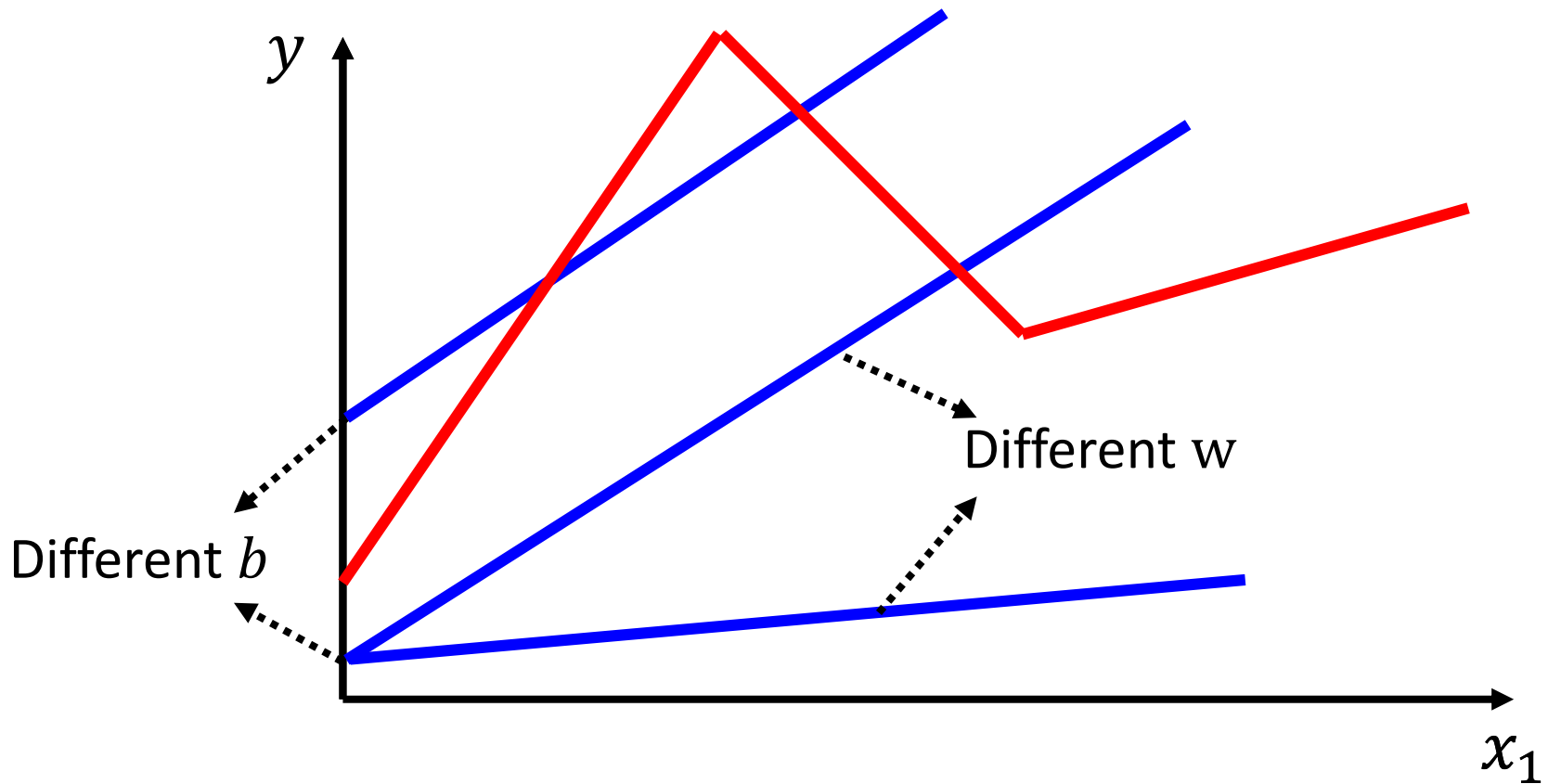
$$L = 0.32k$$

2021

$$L' = 0.46k$$


**Linear models**

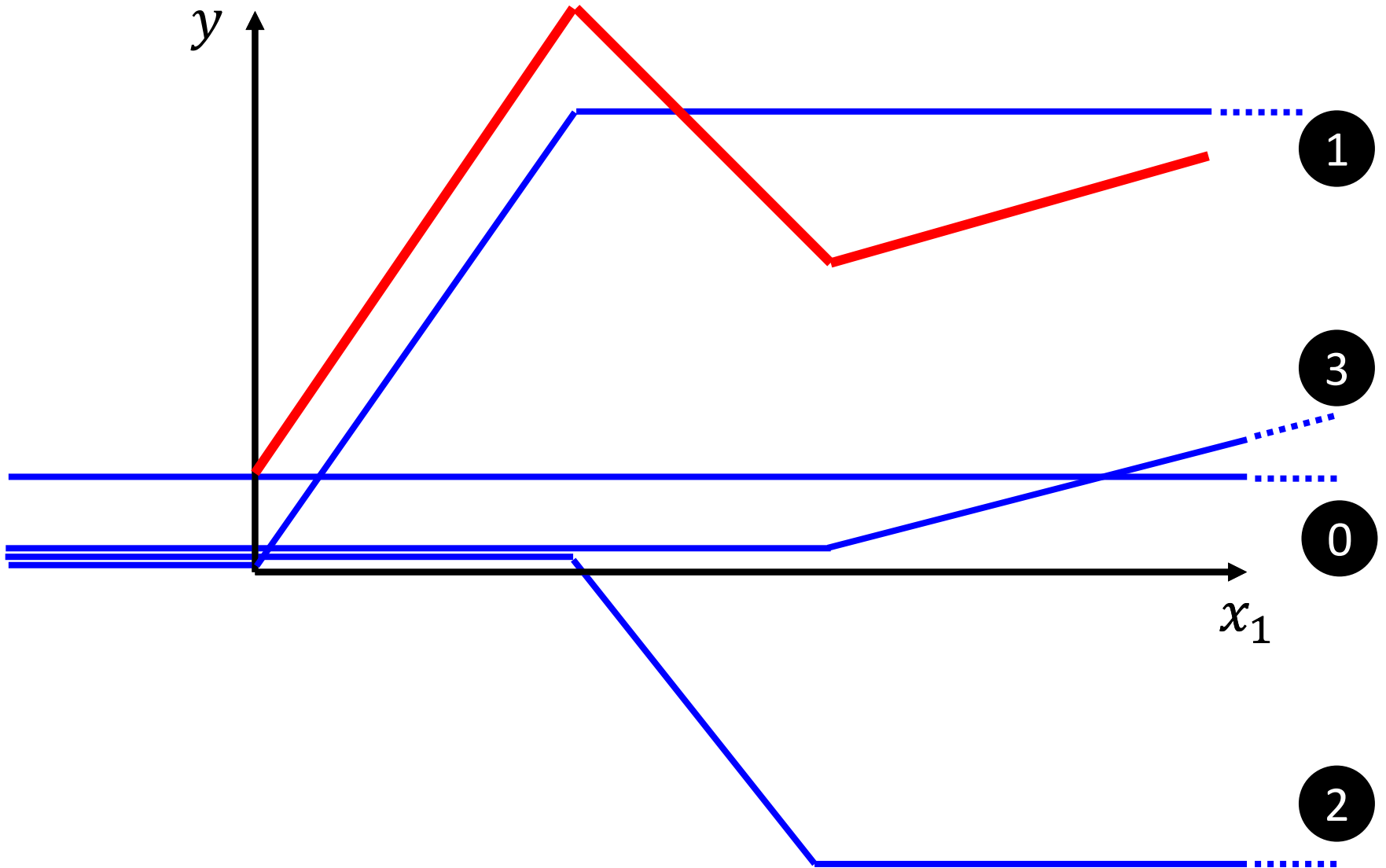
Linear models are too simple ... we need more sophisticated modes.



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

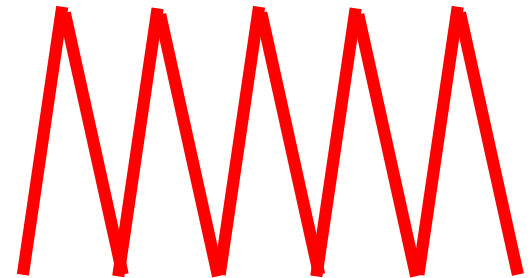
We need a more flexible model!

red curve = constant + sum of a set of 



# All Piecewise Linear Curves

= constant + sum of a set of

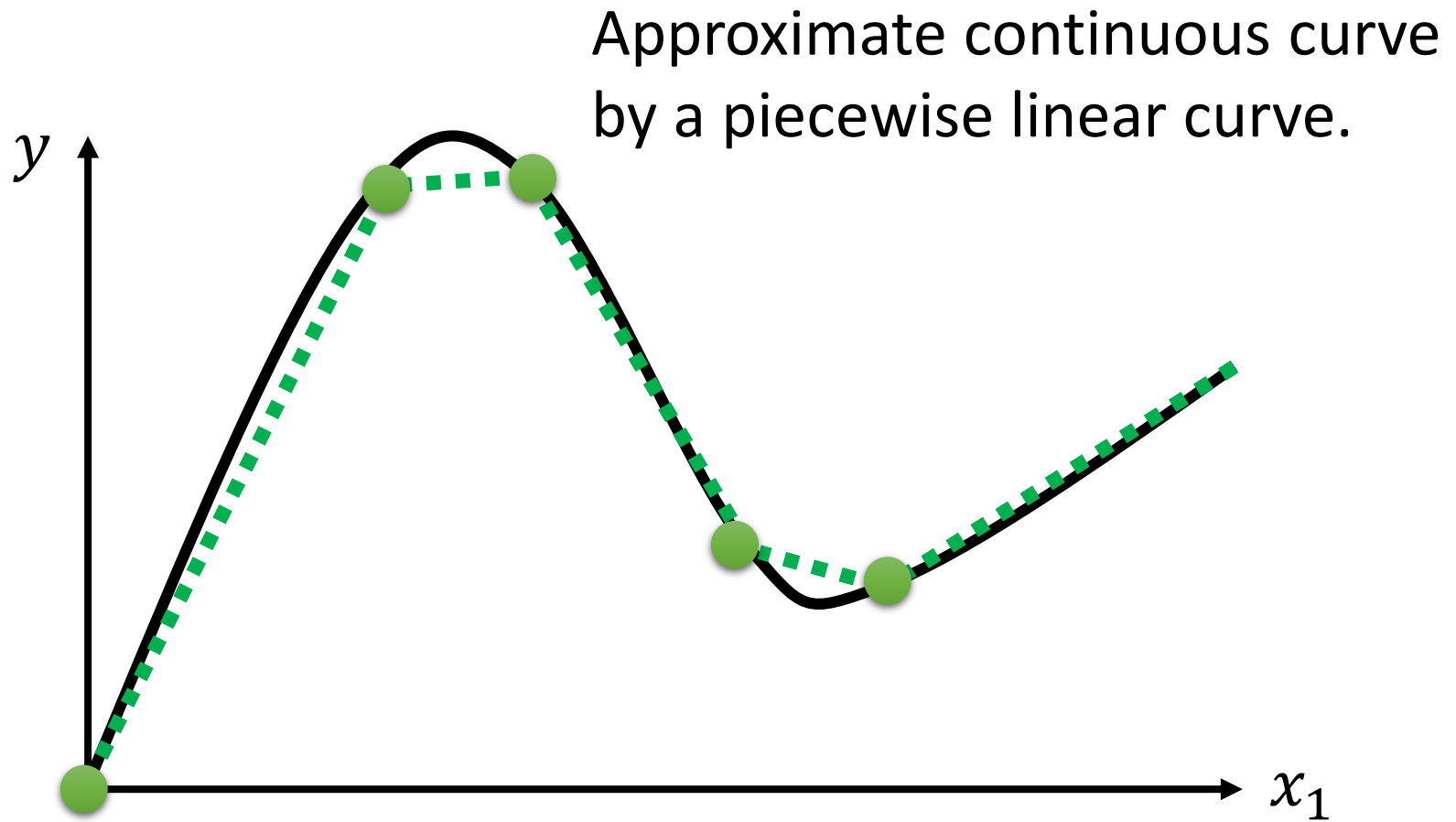


More pieces require more






# Beyond Piecewise Linear?

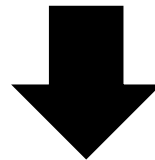


To have good approximation, we need sufficient pieces.

red curve = constant + sum of a set of 

How to represent  
this function?

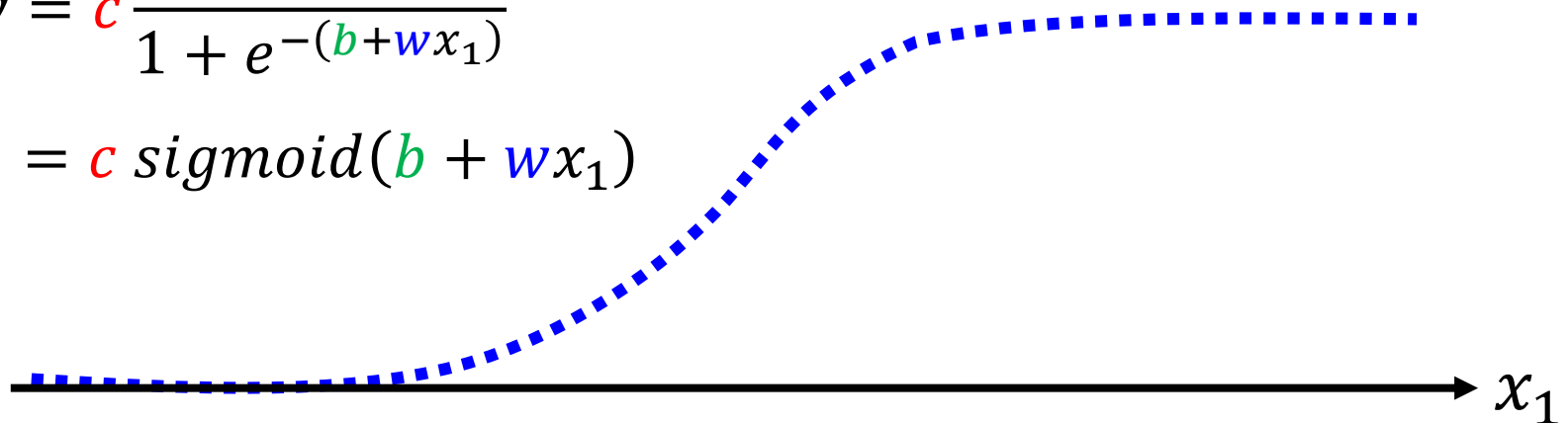
Hard Sigmoid

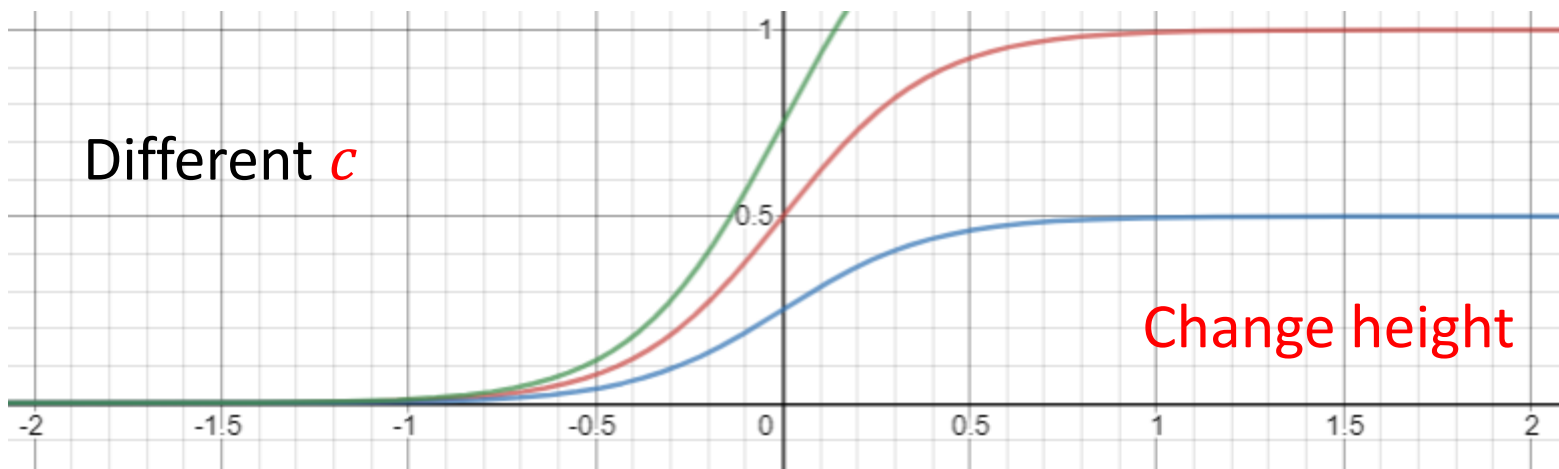
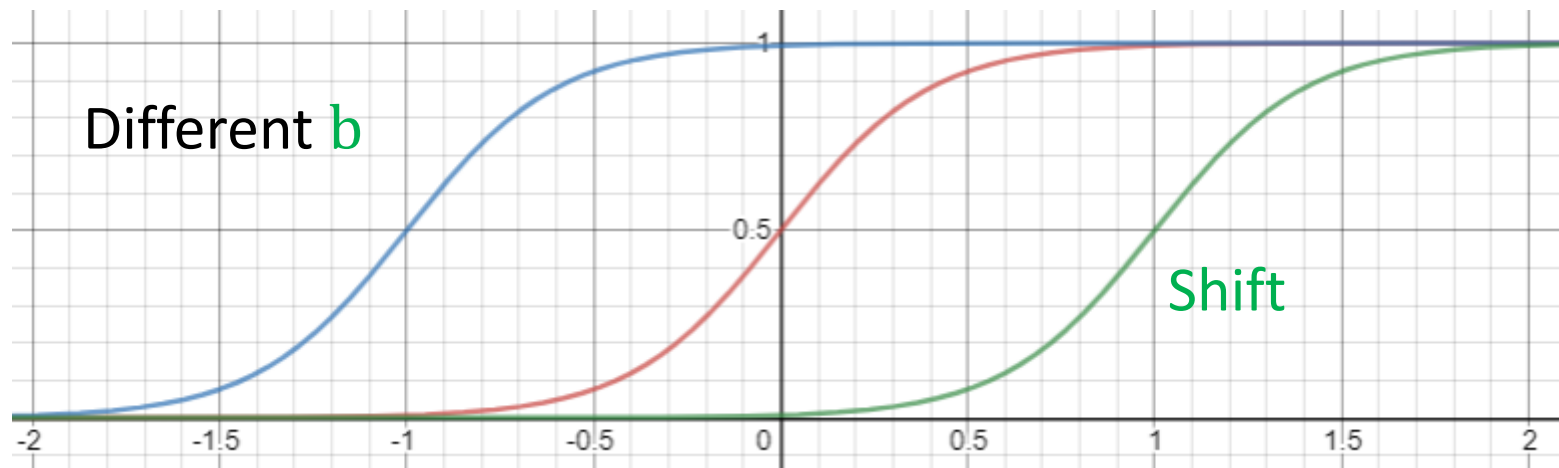
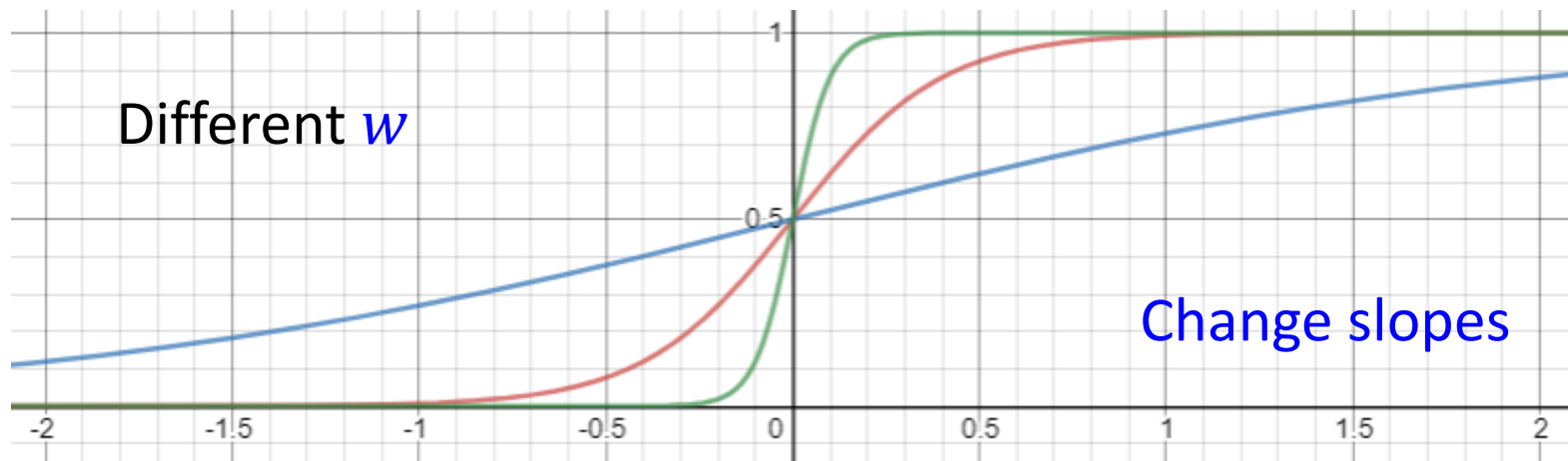


Sigmoid Function

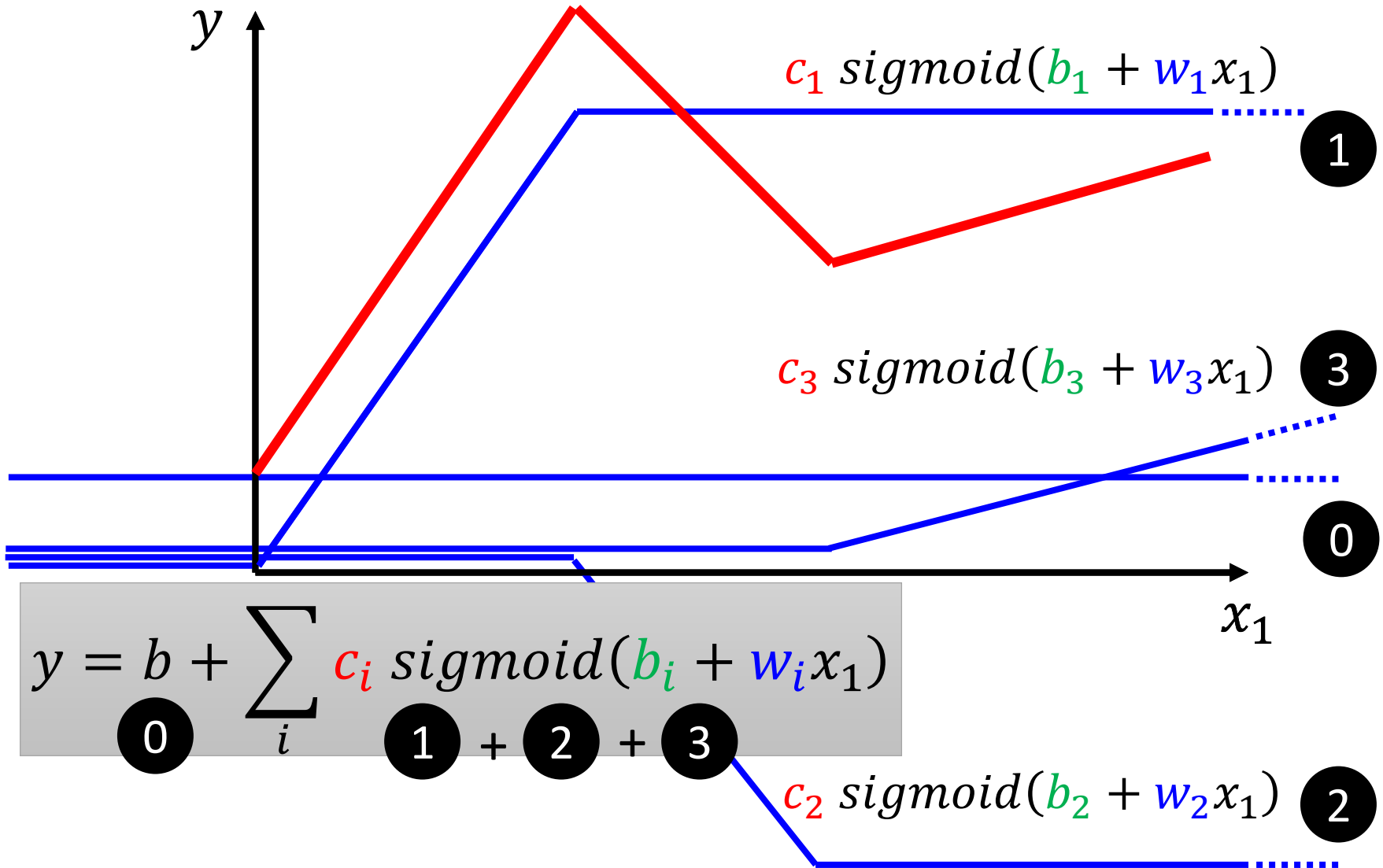
$$y = c \frac{1}{1 + e^{-(b + wx_1)}}$$

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






red curve = sum of a set of  + constant




# New Model: More Features

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

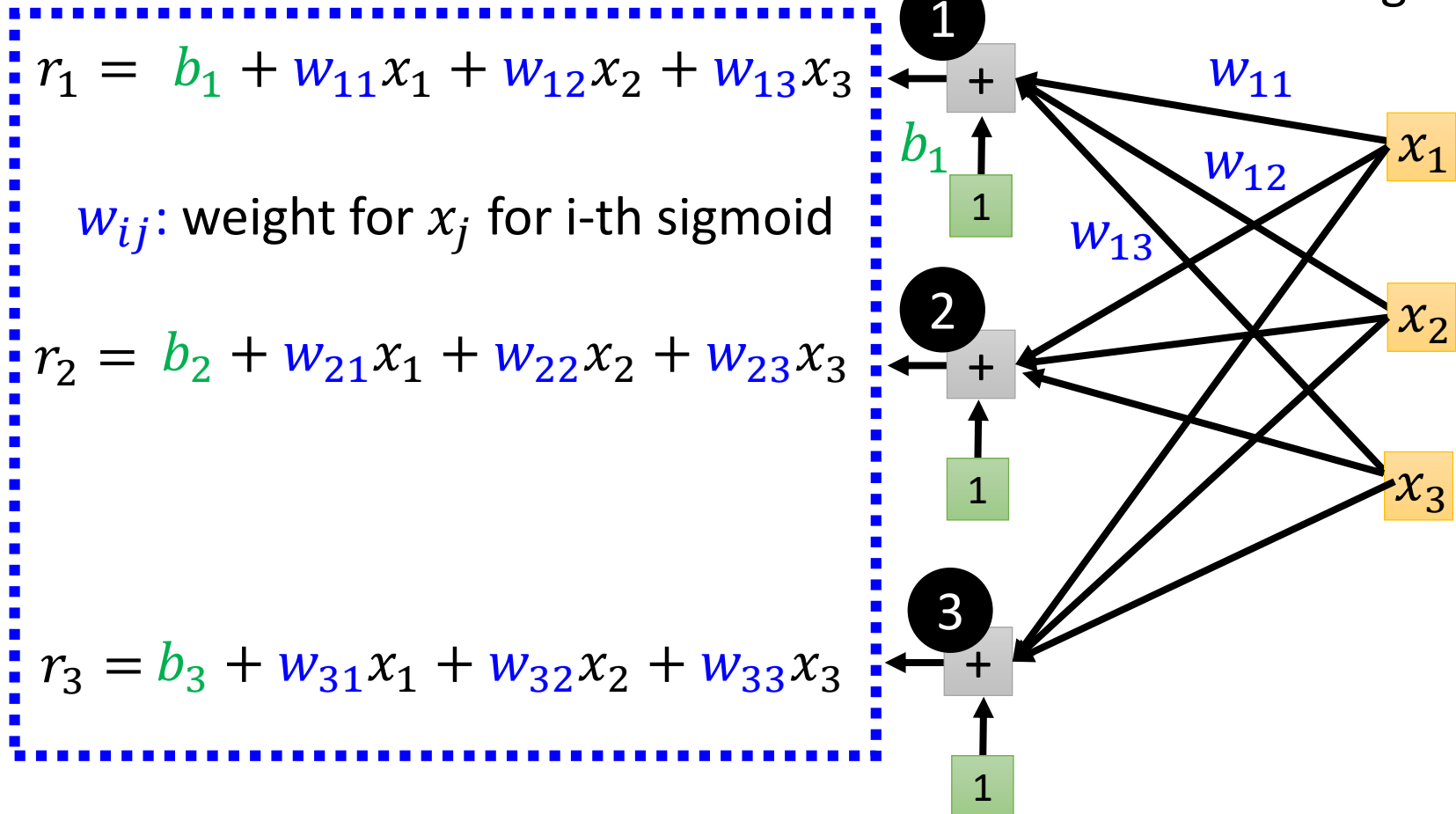

$$y = b + \sum_i \textcolor{red}{c_i} \textit{sigmoid}(\underline{\textcolor{green}{b_i} + \textcolor{blue}{w_i}x_1})$$

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


$$y = b + \sum_i \textcolor{red}{c_i} \textit{sigmoid} \left( \underline{\textcolor{green}{b_i} + \sum_j \textcolor{blue}{w_{ij}} x_j} \right)$$

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

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



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

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

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

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

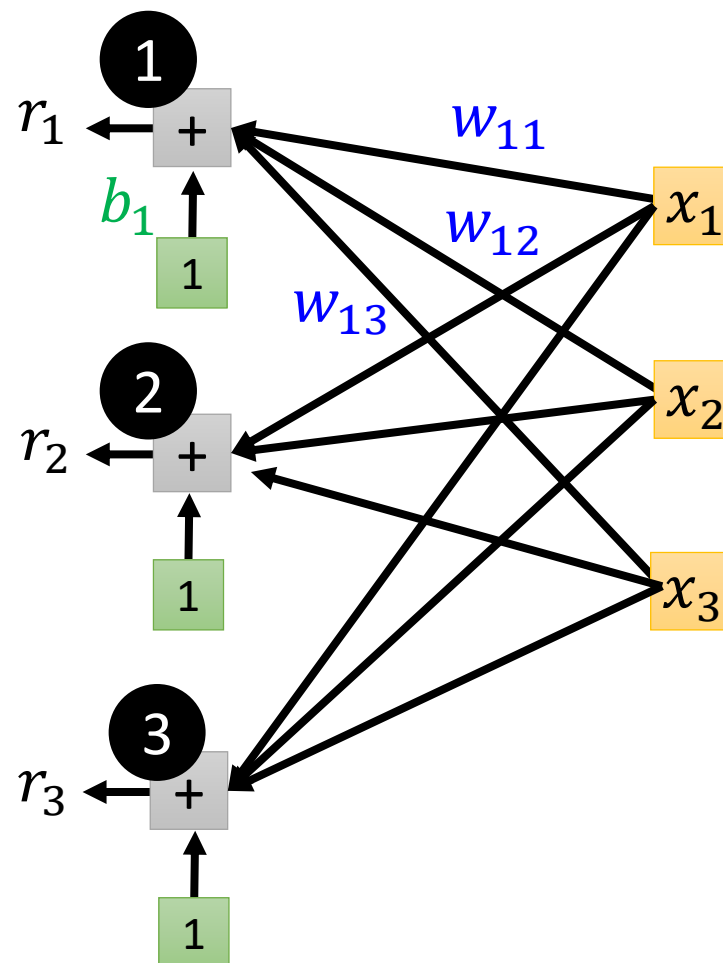
$$\begin{bmatrix} r_1 \\ r_2 \\ r_3 \end{bmatrix} = \begin{bmatrix} \textcolor{green}{b}_1 \\ \textcolor{green}{b}_2 \\ \textcolor{green}{b}_3 \end{bmatrix} + \begin{bmatrix} \textcolor{blue}{w}_{11} & \textcolor{blue}{w}_{12} & \textcolor{blue}{w}_{13} \\ \textcolor{blue}{w}_{21} & \textcolor{blue}{w}_{22}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

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

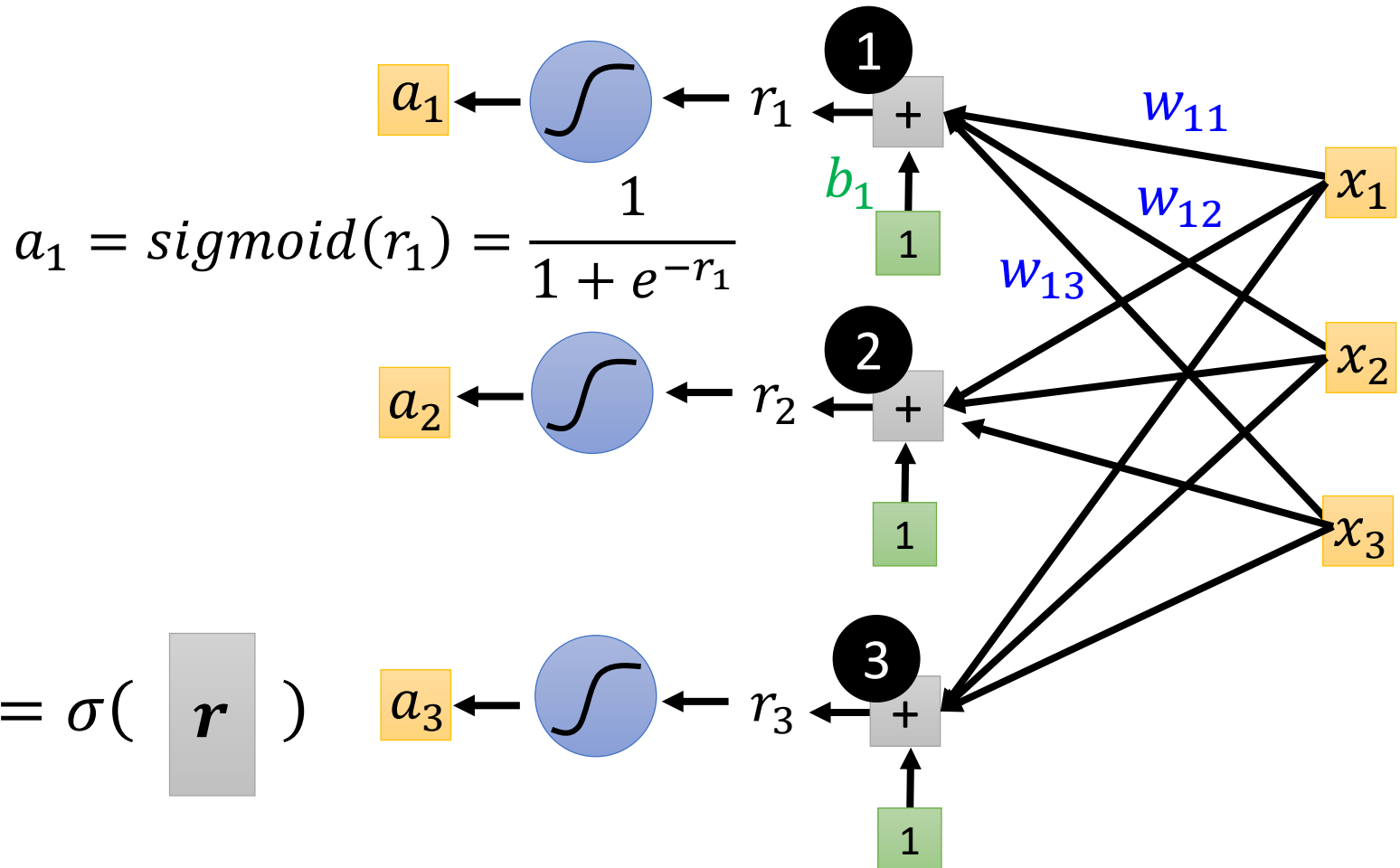


$$y = b + \sum_i \textcolor{red}{c}_i \textit{sigmoid} \left( \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right) \quad \begin{array}{l} i: 1,2,3 \\ j: 1,2,3 \end{array}$$

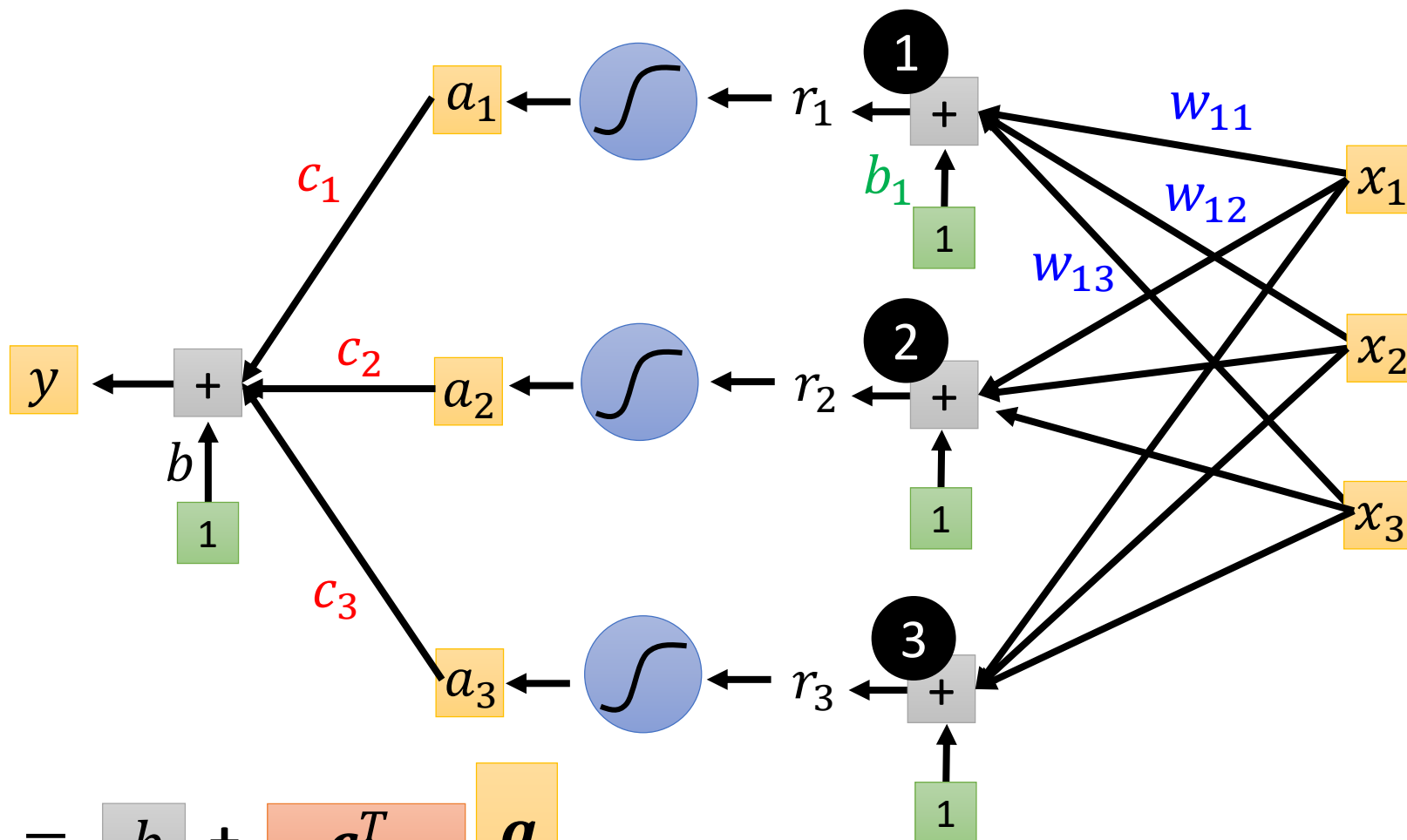
$$\mathbf{r} = \mathbf{b} + \mathbf{W} \mathbf{x}$$



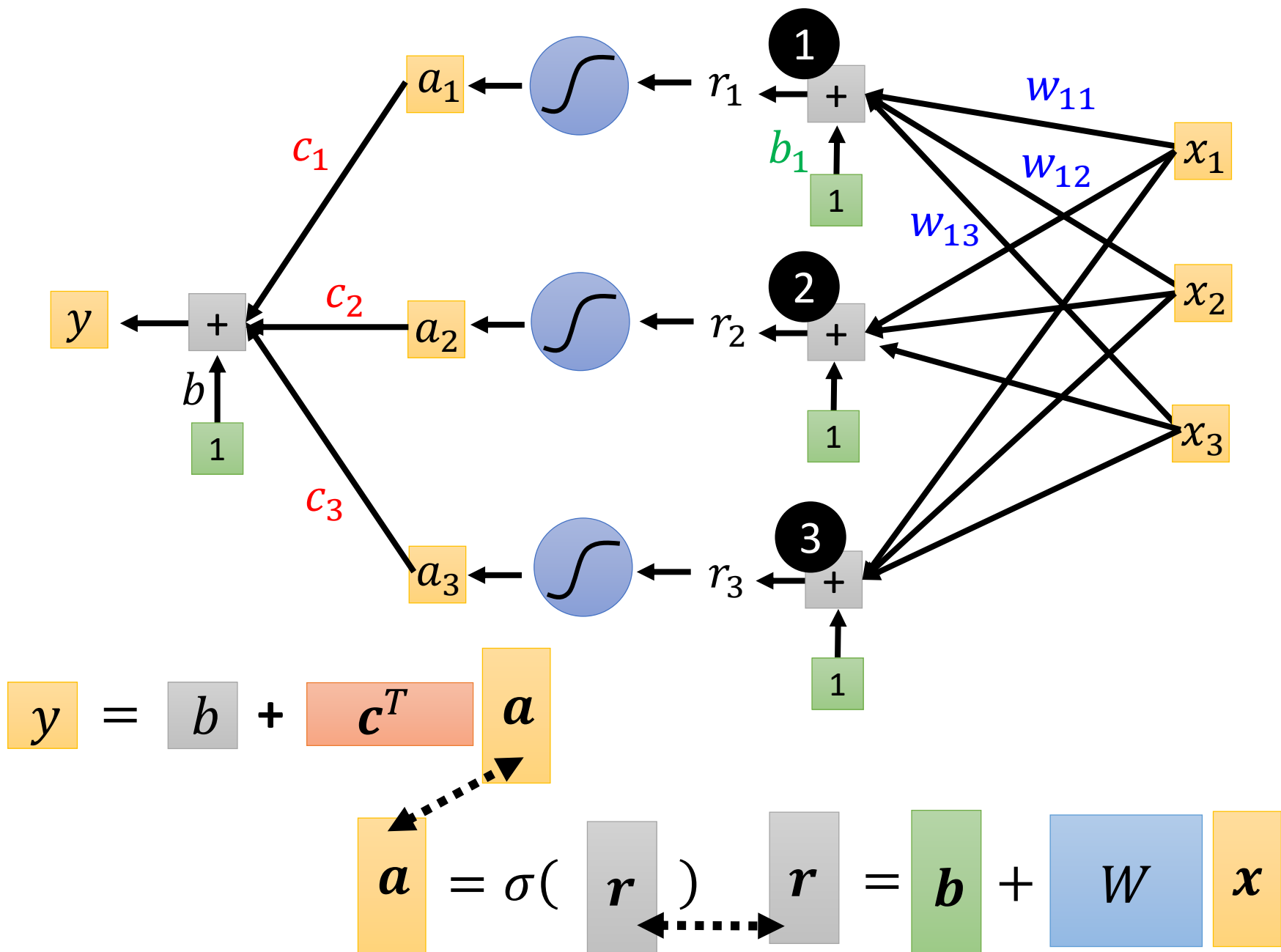
$$y = b + \sum_i \textcolor{red}{c}_i \textcolor{blue}{sigmoid} \left( \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right) \quad \begin{array}{l} i: 1,2,3 \\ j: 1,2,3 \end{array}$$

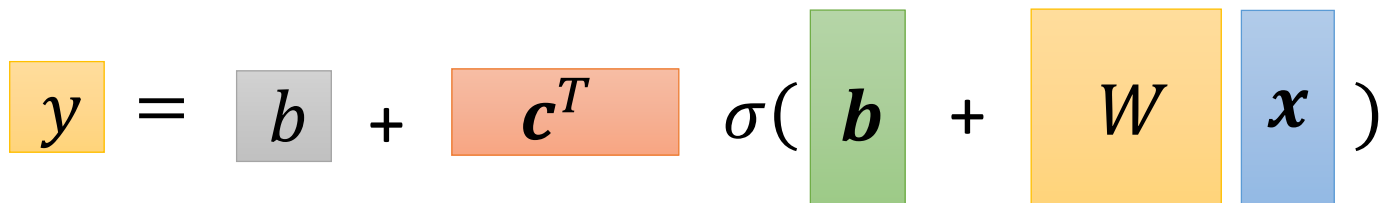


$$y = b + \sum_i \mathbf{c}_i \operatorname{sigmoid} \left( \mathbf{b}_i + \sum_j \mathbf{w}_{ij} x_j \right) \quad \begin{array}{l} i: 1,2,3 \\ j: 1,2,3 \end{array}$$



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





# Function with unknown parameters

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

$x$  feature

Unknown parameters

$W$

$b$

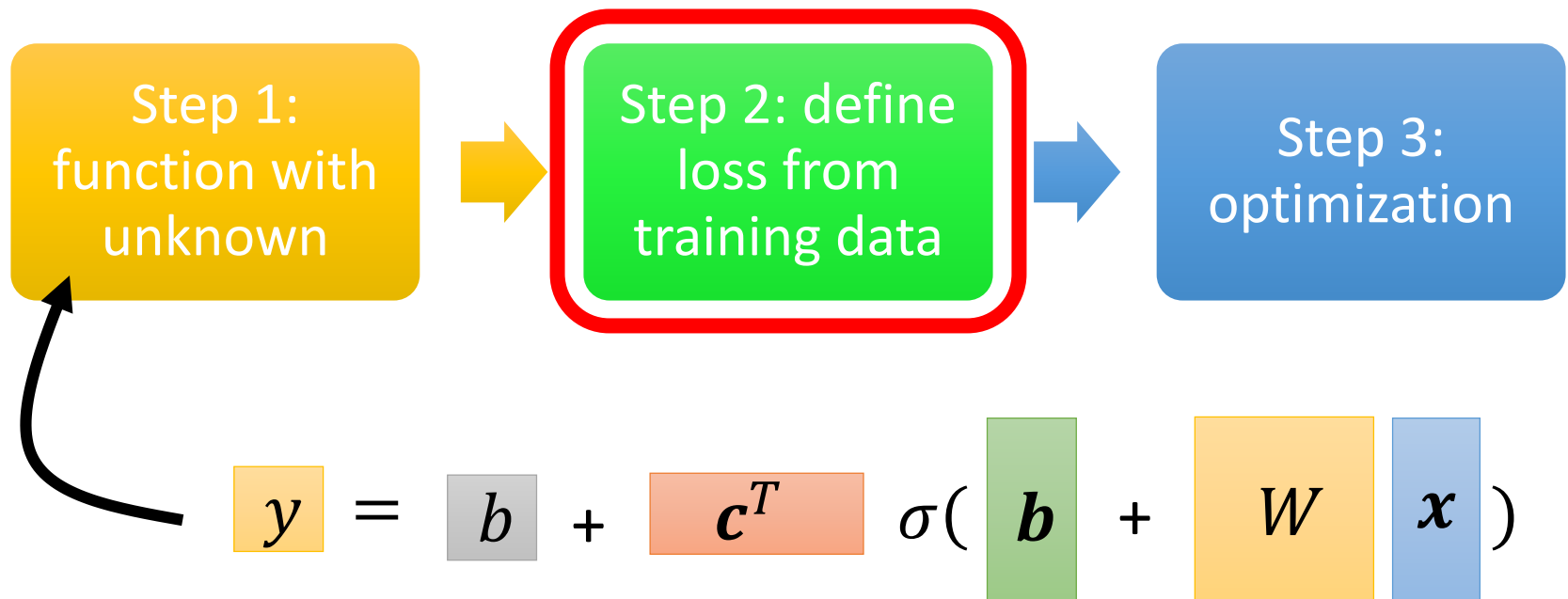
$c^T$

$b$

Rows of  $W$

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

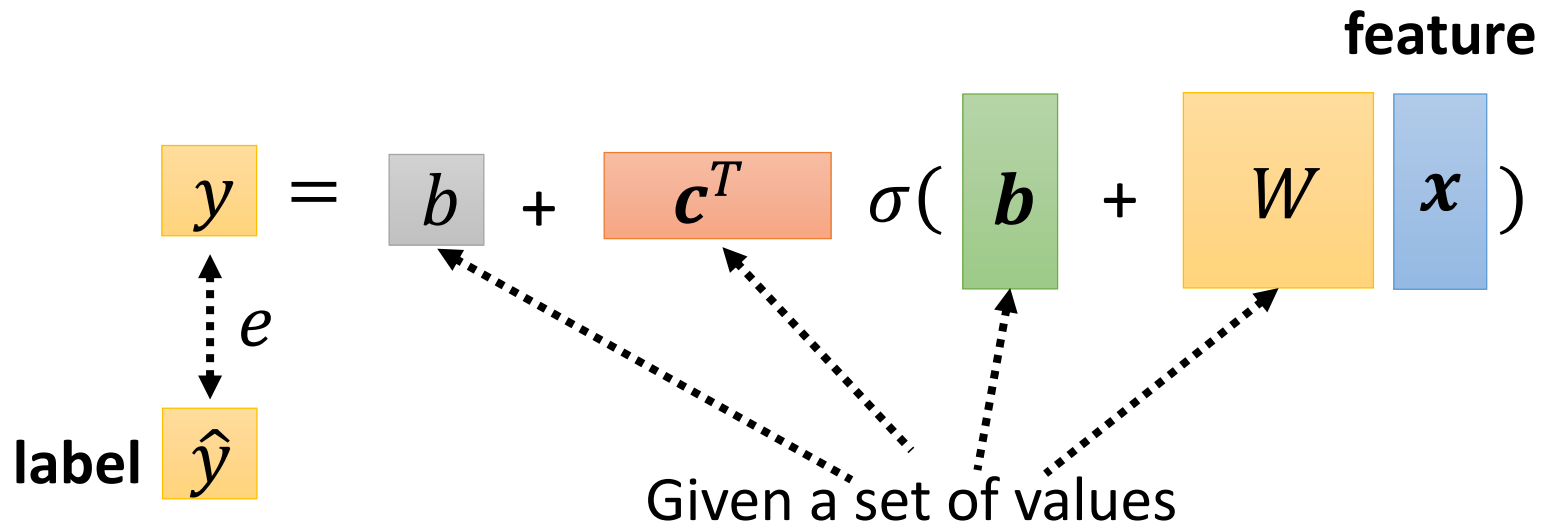
# Back to ML Framework





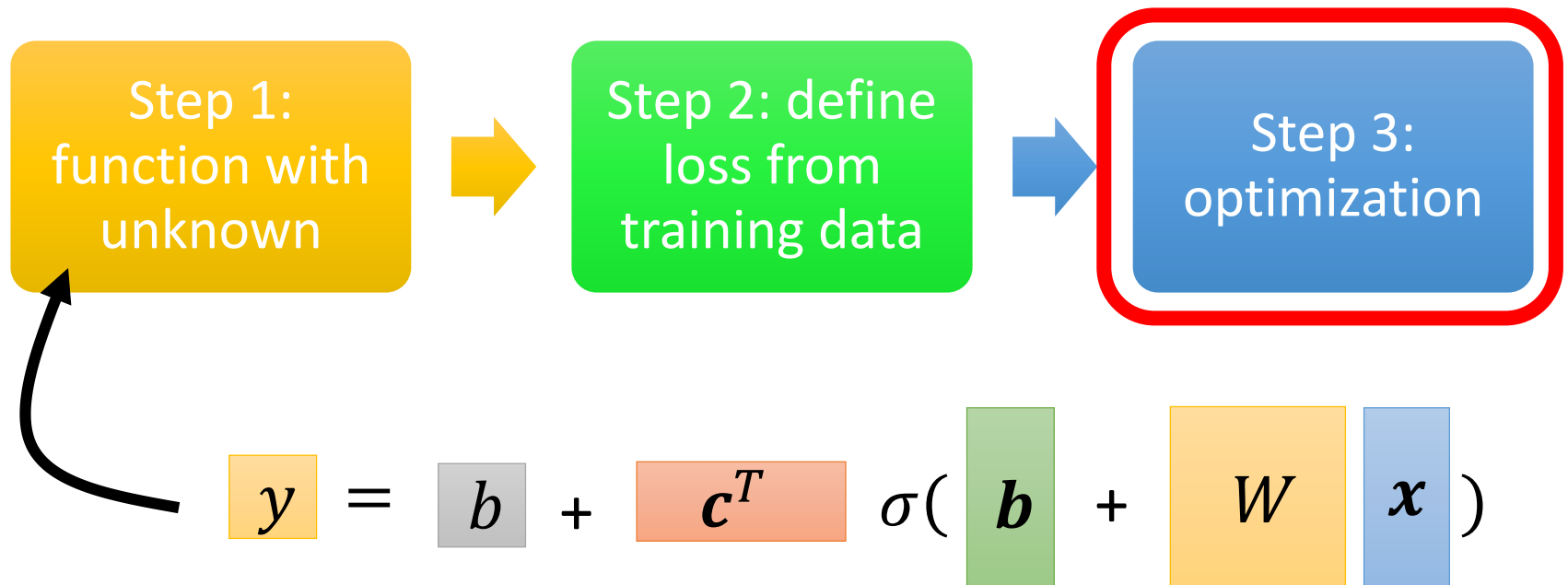
# LOSS

- 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

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L$$

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

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

$$\text{gradient } \mathbf{g} = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \frac{\partial L}{\partial \theta_2} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix} \quad \begin{bmatrix} \theta_1^1 \\ \theta_2^1 \\ \vdots \end{bmatrix} \leftarrow \begin{bmatrix} \theta_1^0 \\ \theta_2^0 \\ \vdots \end{bmatrix} - \begin{bmatrix} \eta \frac{\partial L}{\partial \theta_1} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \eta \frac{\partial L}{\partial \theta_2} \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^0} \\ \vdots \end{bmatrix}$$

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

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

# Optimization of New Model

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L$$

➤ (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

$$\theta^* = \arg \min_{\theta} L$$

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

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

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

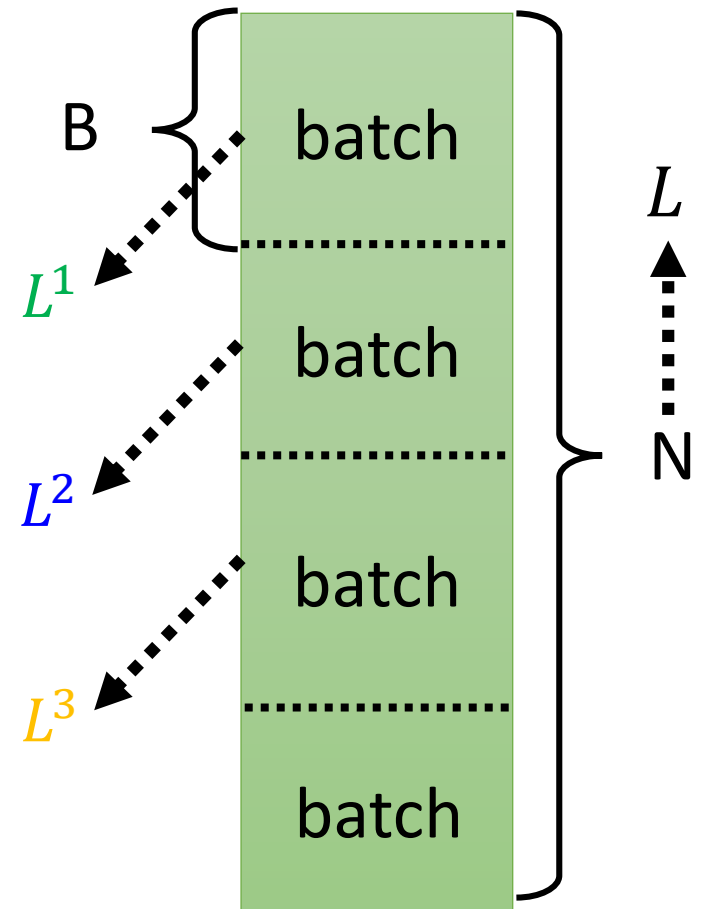
➤ Compute gradient  $\mathbf{g} = \nabla L^2(\theta^1)$

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

➤ 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**?

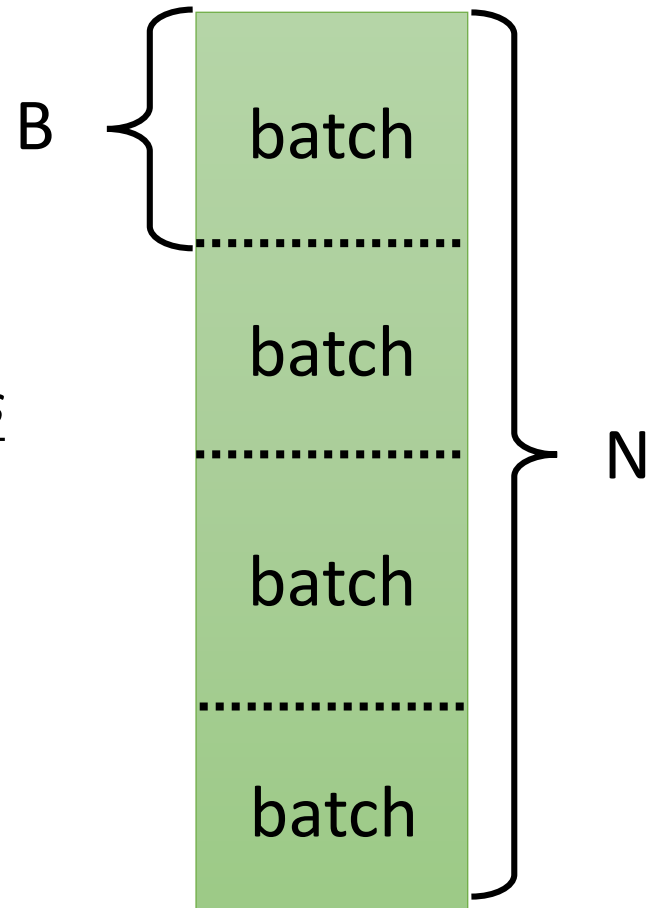
1,000 updates

## Example 2

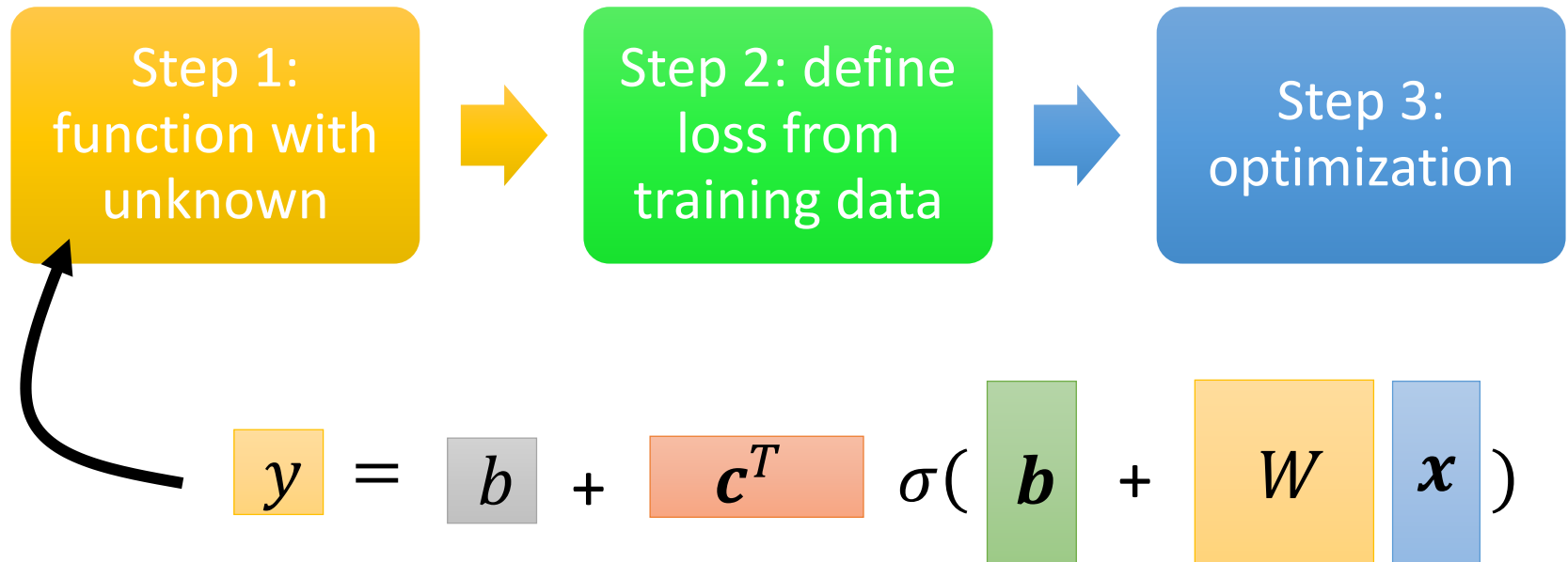
- 1,000 examples ( $N = 1,000$ )
- Batch size is 100 ( $B = 100$ )

How many update in **1 epoch**?

10 updates



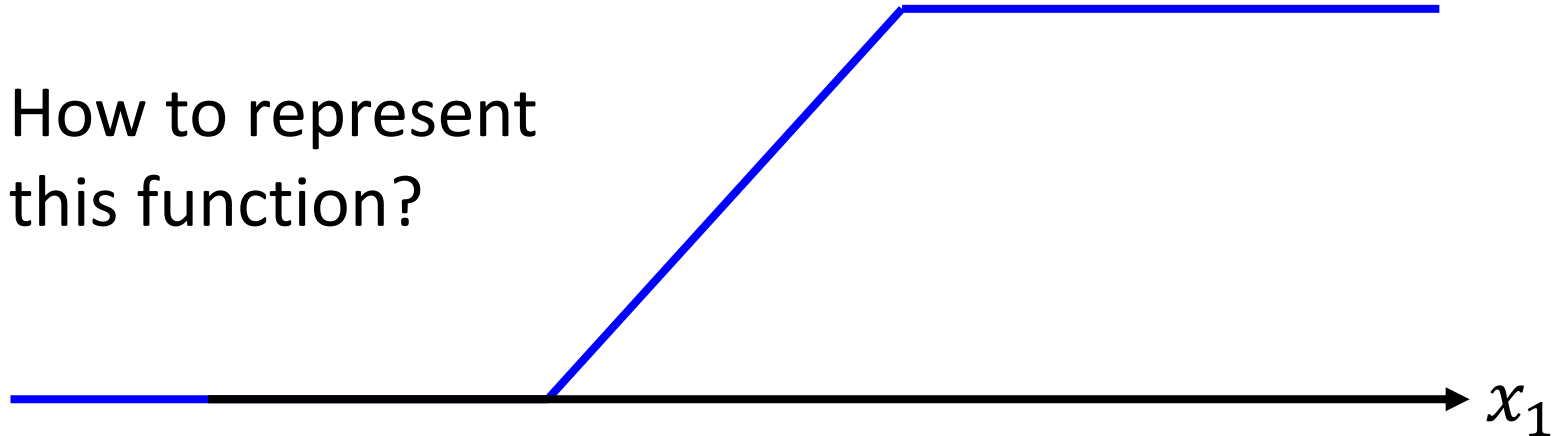
# Back to ML Framework



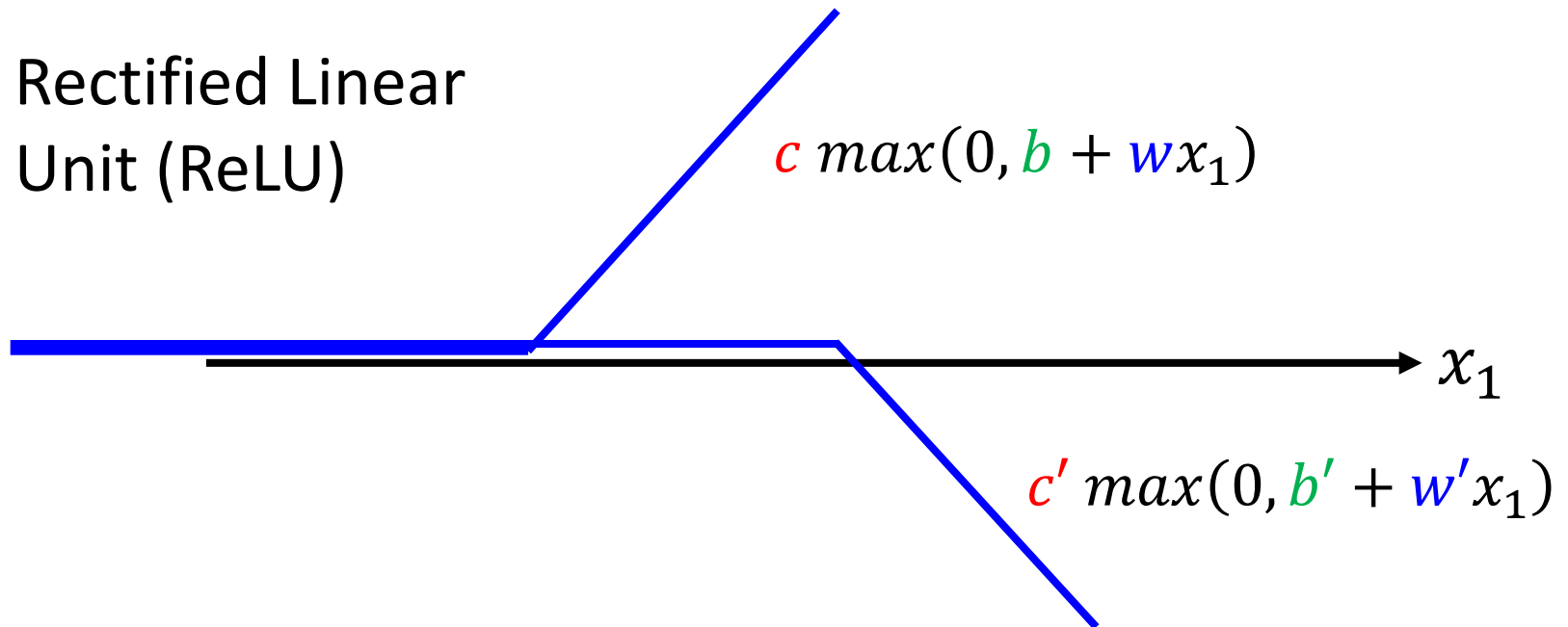
More variety of models ...

# Sigmoid $\rightarrow$ ReLU

How to represent  
this function?




Rectified Linear  
Unit (ReLU)





# Sigmoid $\rightarrow$ ReLU

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


**Activation function**

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

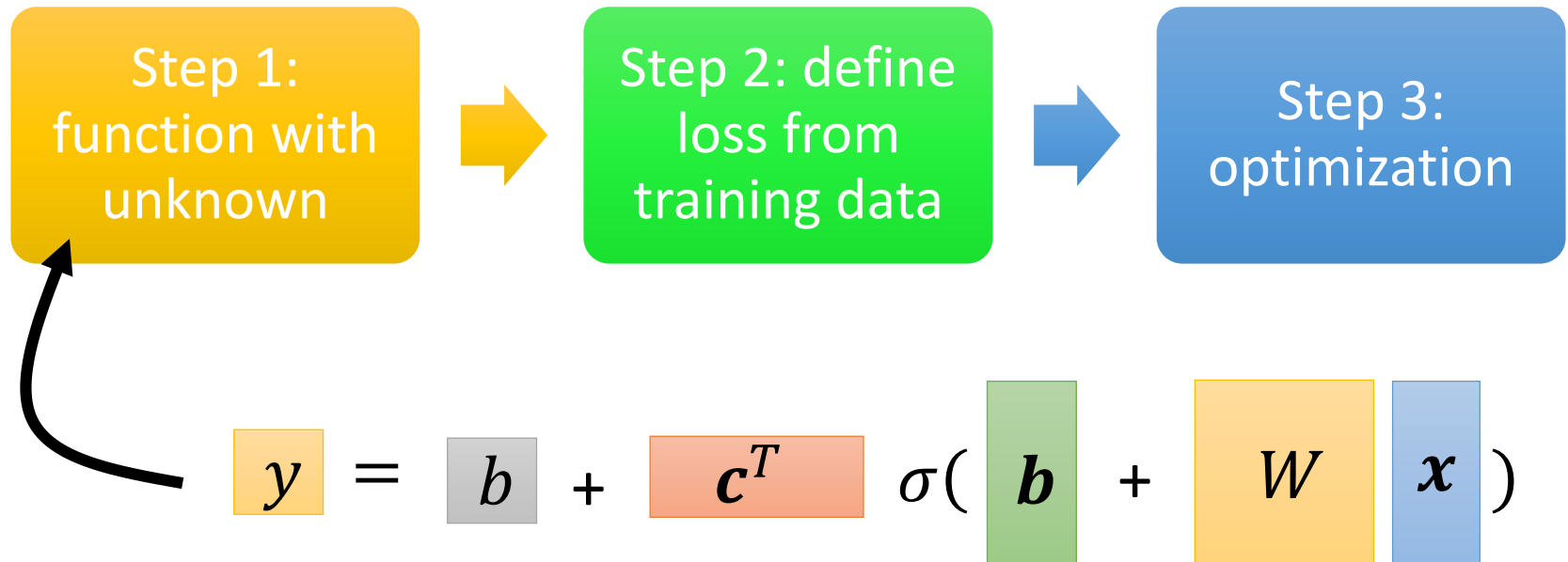
Which one is better?

# Experimental Results

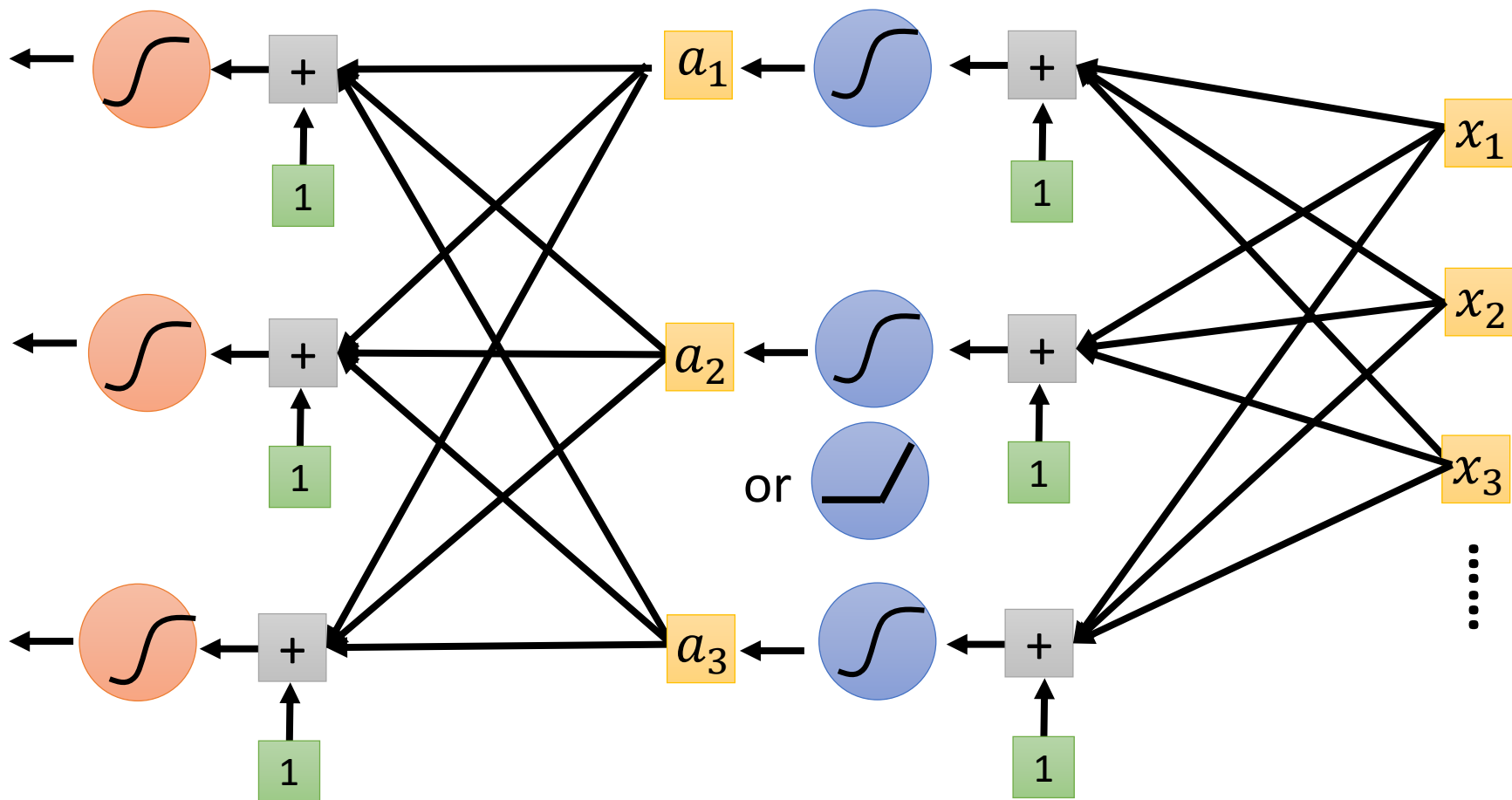
$$y = b + \sum_i \textcolor{red}{c}_i \max \left( 0, \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right)$$

	linear	10 ReLU	100 ReLU	1000 ReLU
2017 – 2020	0.32k	0.32k	0.28k	0.27k
2021	0.46k	0.45k	0.43k	0.43k

# Back to ML Framework



Even more variety of models ...



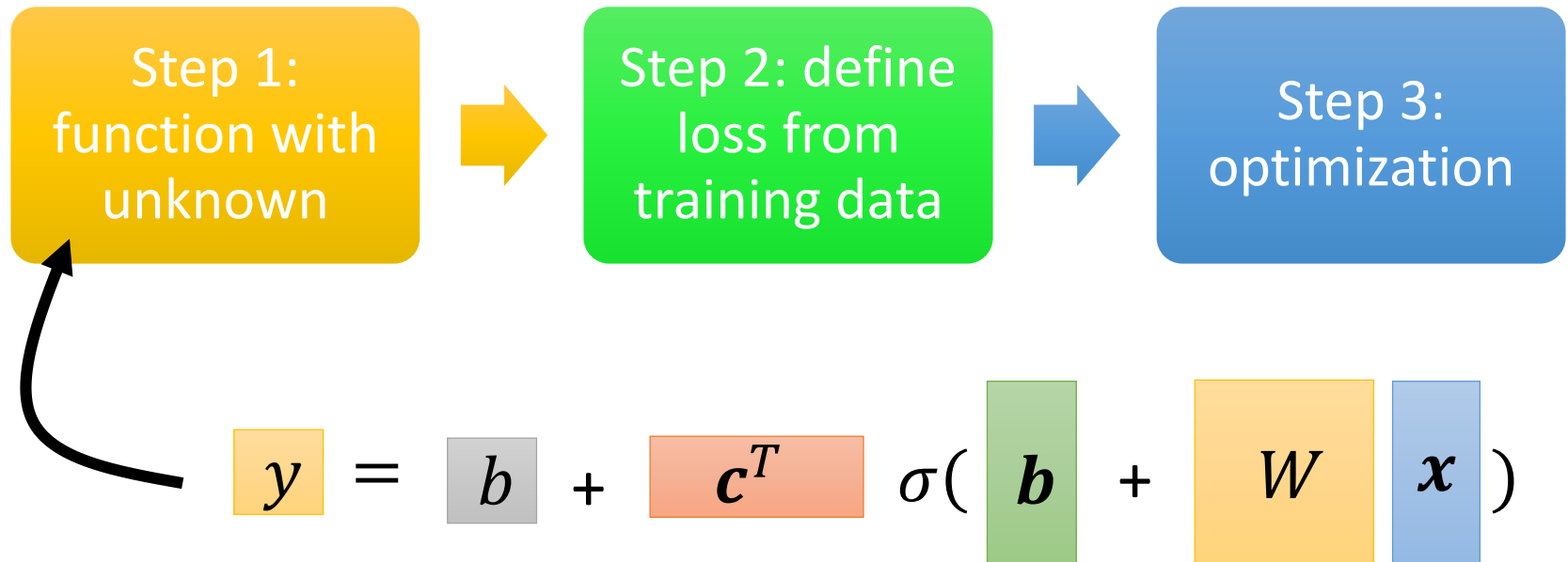
$$\begin{aligned}
 \mathbf{a}' &= \sigma( \mathbf{b}' + \mathbf{W}' \mathbf{a} ) & \mathbf{a} &= \sigma( \mathbf{b} + \mathbf{W} \mathbf{x} )
 \end{aligned}$$

# Experimental Results

- Loss for multiple hidden layers
  - 100 ReLU for each layer
  - input features are the no. of views in the past 56 days

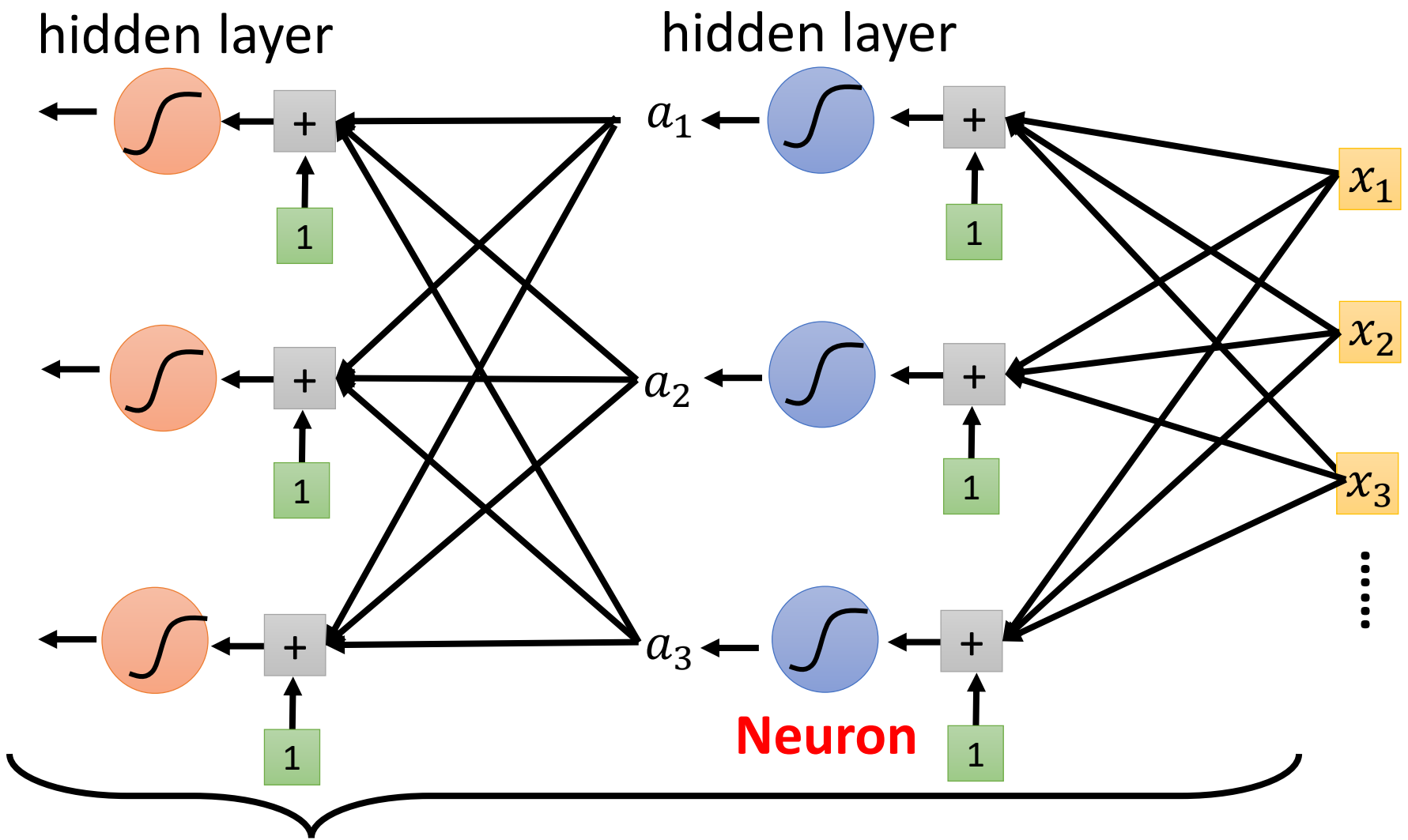
	1 layer	2 layer	3 layer	4 layer
2017 – 2020	0.28k	0.18k	0.14k	0.10k
2021	0.43k	0.39k	0.38k	0.44k

# Back to ML Framework



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

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



**Neural Network**

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

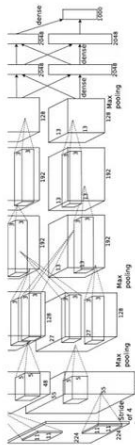
Many layers means **Deep** ➡ Deep Learning

# Deep = Many hidden layers

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

8 layers

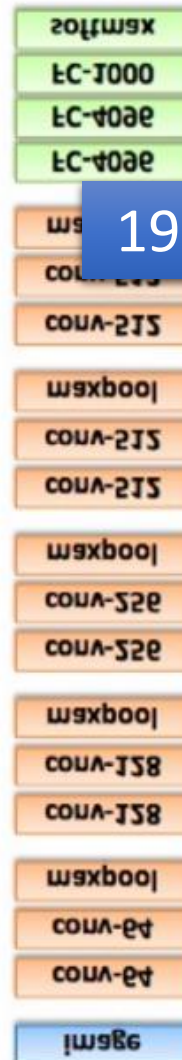
16.4%



AlexNet (2012)

19 layers

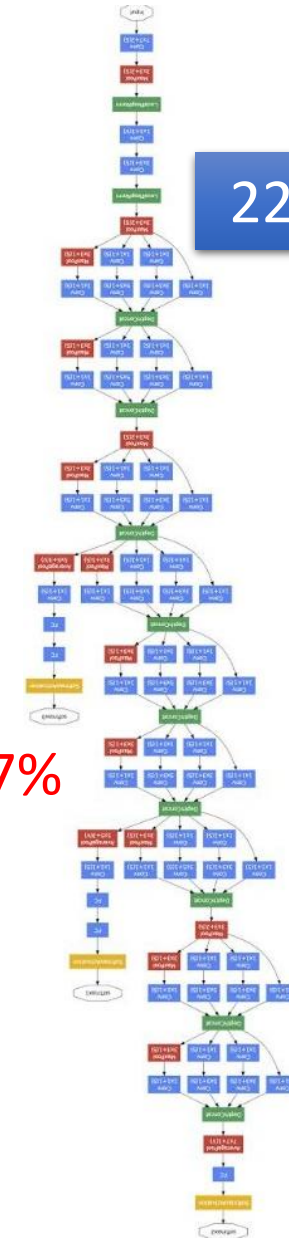
7.3%



VGG (2014)

22 layers

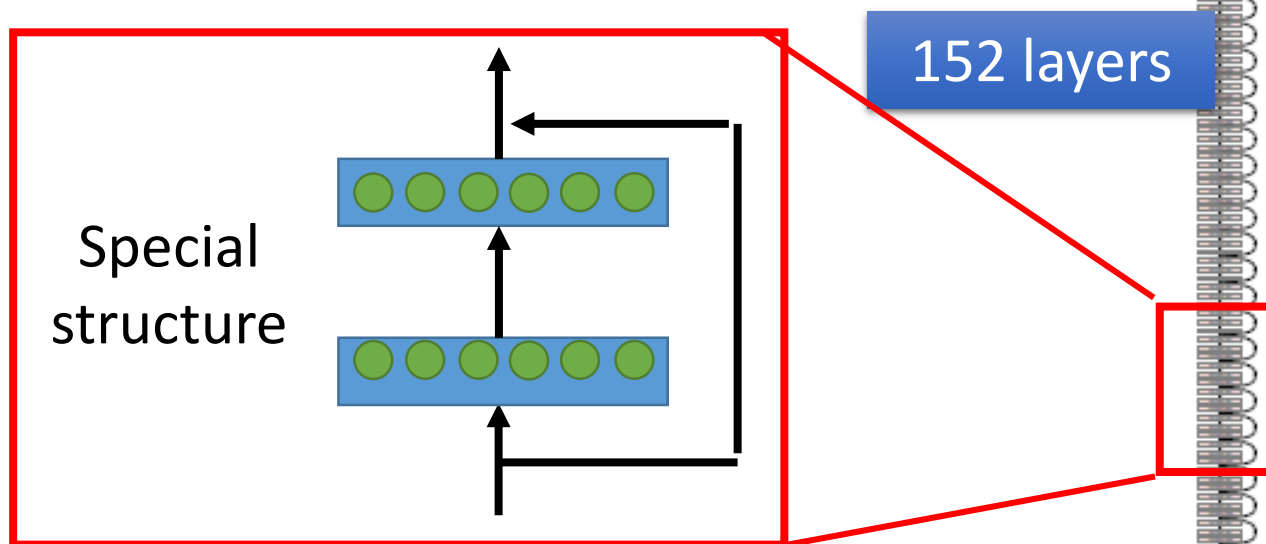
6.7%



GoogleNet (2014)



# Deep = Many hidden layers



152 layers

101 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)

Taipei  
101



# Why don't we go deeper?

- Loss for multiple hidden layers
  - 100 ReLU for each layer
  - input features are the no. of views in the past 56 days

	1 layer	2 layer	3 layer
2017 – 2020	0.28k	0.18k	0.14k
2021	0.43k	0.39k	0.38k

# Why don't we go deeper?

- Loss for multiple hidden layers
  - 100 ReLU for each layer
  - input features are the no. of views in the past 56 days

	1 layer	2 layer	3 layer	4 layer
2017 – 2020	0.28k	0.18k	0.14k	0.10k
2021	0.43k	0.39k	0.38k	0.44k

Better on training data, worse on unseen data

 **Overfitting**