



北京航空航天大学
COLLEGE OF SOFTWARE 软件学院
BEIHANG UNIVERSITY

人工智能

第6讲：深度学习I

梯度下降与前馈神经网络

张晶

2025年春季

- 参考资料：吴飞，《人工智能导论：模型与算法》，高等教育出版社
- 在线课程：<https://www.icourse163.org/course/ZJU-1003377027?from=searchPage>
- 本部分参考：李宏毅，《机器学习》课程，台湾大学



北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

提纲

一、线性回归与梯度下降

二、前馈神经网络

三、卷积神经网络

四、序列数据模型

五、深度学习应用



北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

提纲

一、线性回归与梯度下降

二、前馈神经网络

三、卷积神经网络

四、序列数据模型

五、深度学习应用

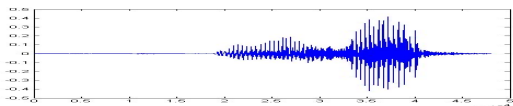


北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

机器学习：拟合函数 $f(x) = y$

- 语音识别

$$f(\text{语音波形}) = \text{“你好”}$$



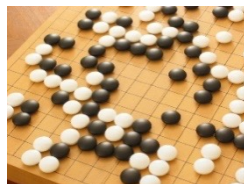
- 图像分类

$$f(\text{猫的图片}) = \text{“猫”}$$



- 围棋游戏

$$f(\text{围棋棋盘}) = \text{“5-5”}$$



(下一步落子位置)

- 从原始数据中提取特征
- 学习映射函数 f
- 通过映射函数 f 将原始数据映射到语义任务空间，即寻找数据和任务目标之间的关系

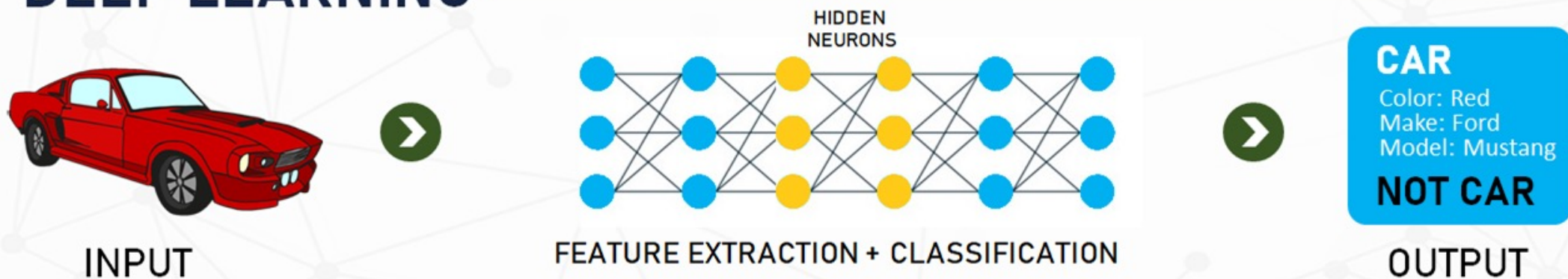


机器学习 v. s. 深度学习

MACHINE LEARNING



DEEP LEARNING





回顾：线性回归

- 任务：网络视频观看量预测
- 训练数据 \mathcal{D} ：2017/01/01 – 2020/12/31

$$\mathcal{D} = \{(x^1, \hat{y}^1), (x^2, \hat{y}^2), \dots, (x^N, \hat{y}^N)\}$$

Diagram illustrating the training data \mathcal{D} structure:

- For (x^1, \hat{y}^1) : x^1 is 2017/01/01 (观看量), \hat{y}^1 is 2017/01/02 (观看量).
- For (x^2, \hat{y}^2) : x^2 is 2017/01/02 (观看量), \hat{y}^2 is 2017/01/03 (观看量).
- For (x^N, \hat{y}^N) : x^N is 2020/12/30 (观看量), \hat{y}^N is 2020/12/31 (观看量).



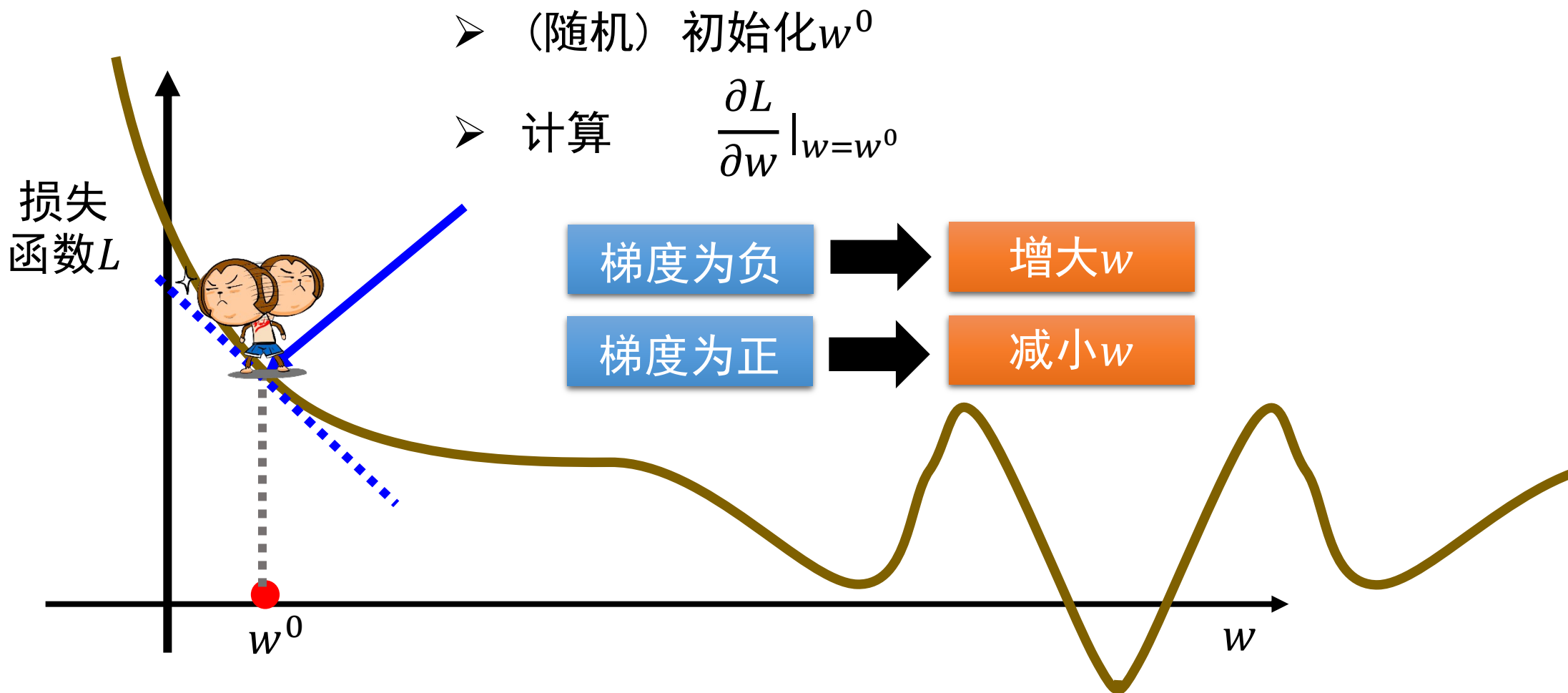
回顾：线性回归

- 训练数据： $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- 学习模型： $f(x_i) = wx_i + b \ (1 \leq i \leq n)$
- 损失函数： $L(w, b) = \sum_{i=1}^n \text{loss}(x_i, y_i) = \sum_{i=1}^n (y_i - (wx_i + b))^2$
- 优化方法： 最小二乘法或梯度下降法



优化方法：梯度下降

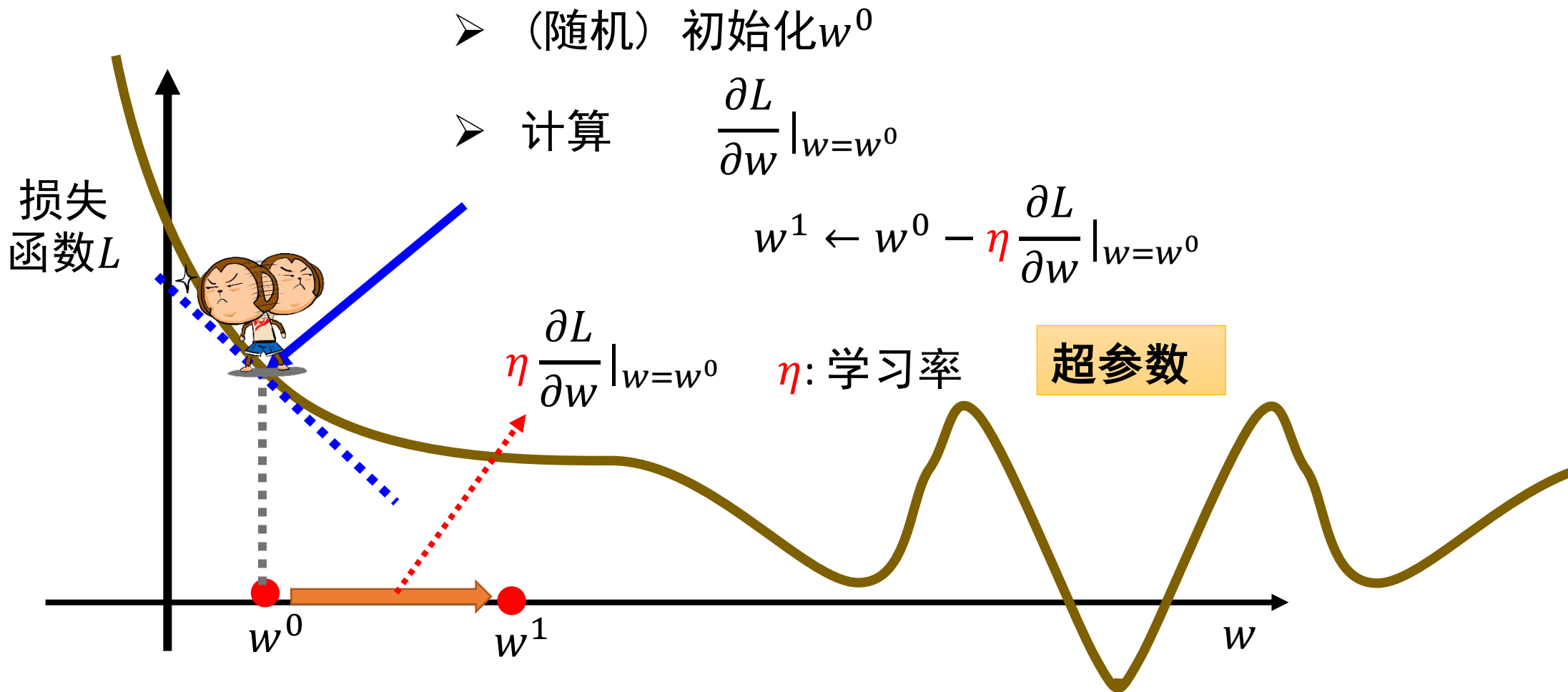
$$w^* = \arg \min_w L$$





优化方法：梯度下降

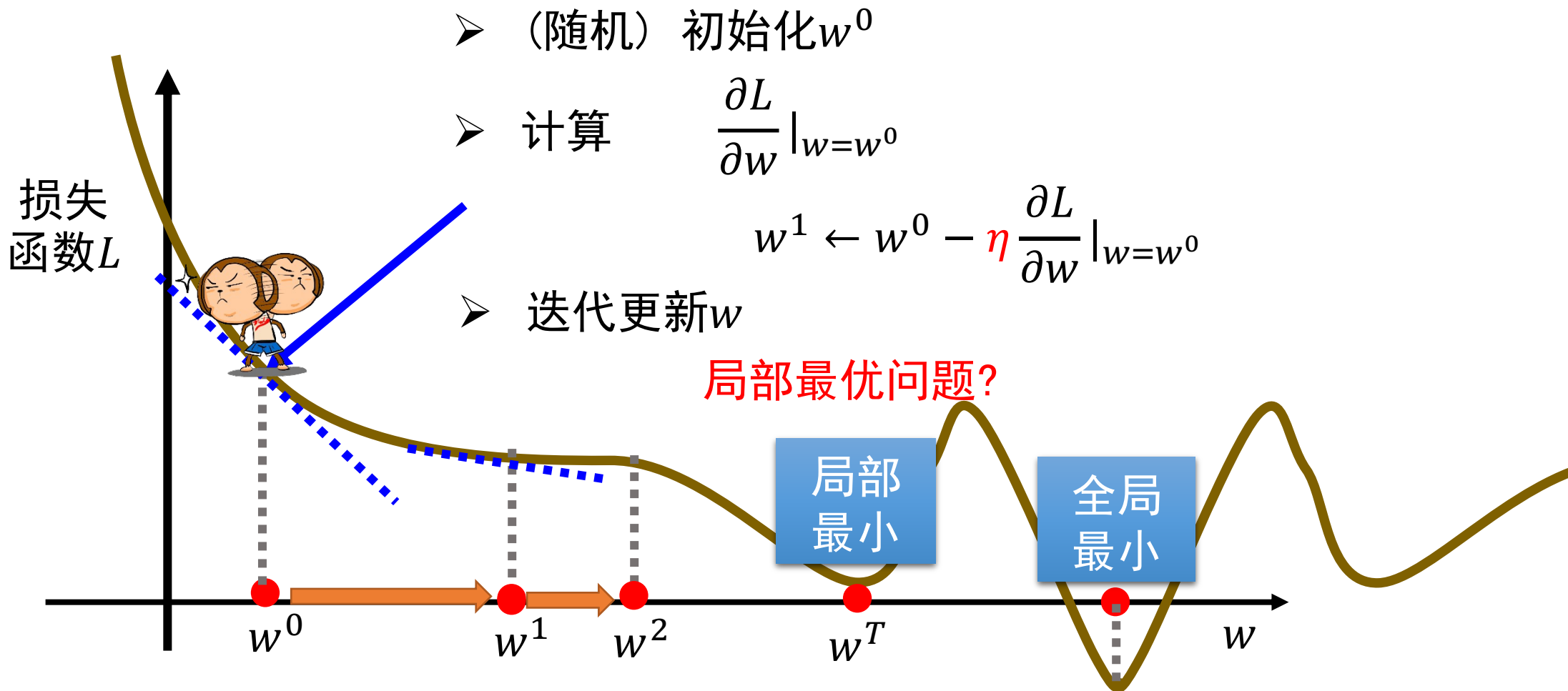
$$w^* = \arg \min_w L$$





优化方法：梯度下降

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





优化方法：梯度下降

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

- (随机) 初始化 w^0, b^0
- 计算

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



深度学习框架下一行代码足以实现！

- 迭代更新 w 和 b

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

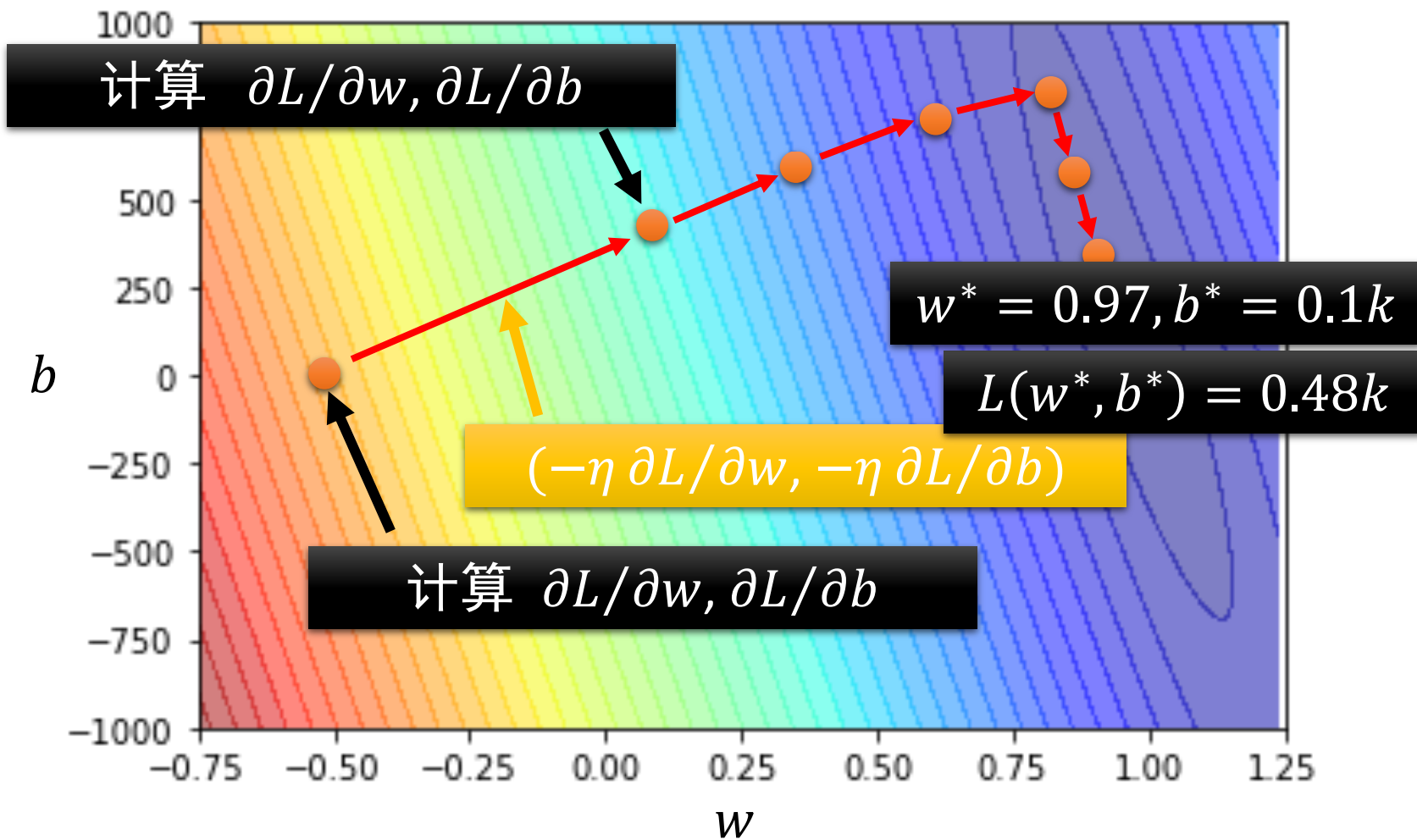
$$b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$$



优化方法：梯度下降

模型 $y = b + wx_1$

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





线性回归：一元线性回归

- 训练数据： $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- 学习模型： $f(x_i) = wx_i + b \ (1 \leq i \leq n)$
- 损失函数： $L(w, b) = \sum_{i=1}^n \text{loss}(x_i, y_i) = \sum_{i=1}^n (y_i - (wx_i + b))^2$
- 优化方法： 梯度下降法得到 $w^* = 0.97, b^* = 0.1k$
- 结果预测： $f(x) = 0.97x + 100$



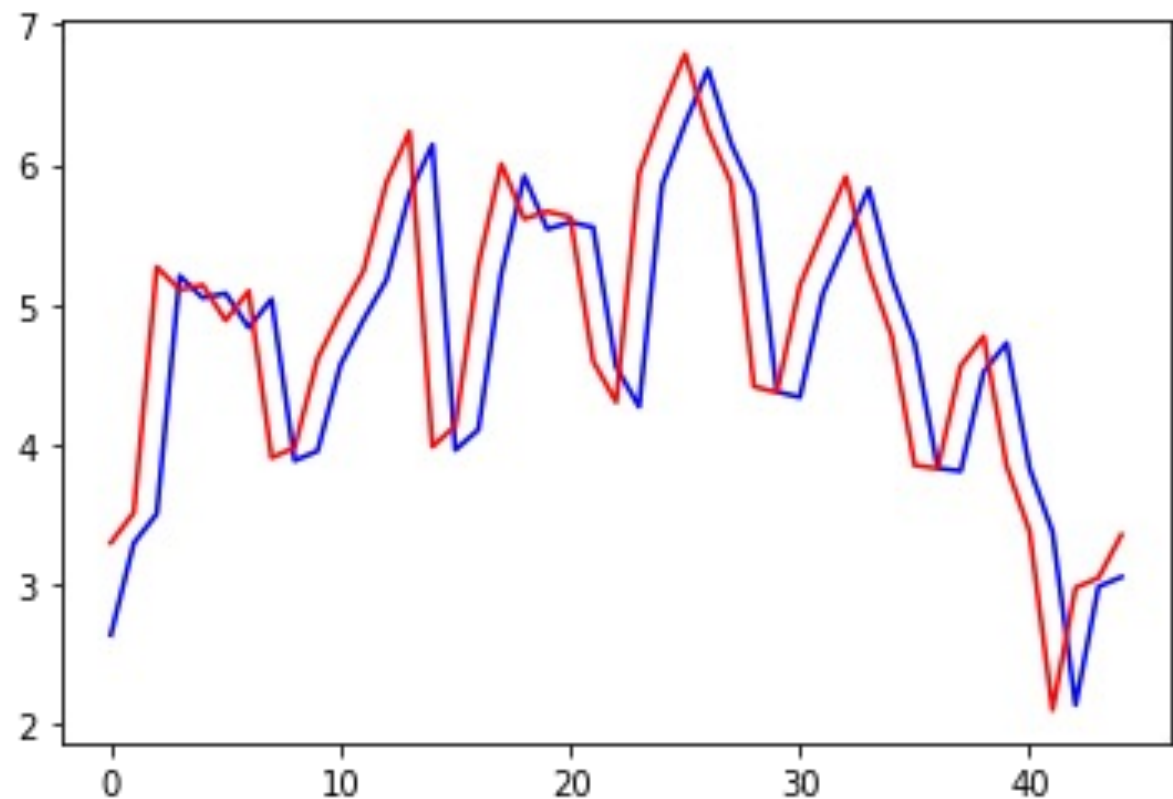
北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

$$y = 0.1k + 0.97x_1$$

Red: 真实观看数量

blue: 预测观看数量

观看数量
(k)



2021/01/01

2021/02/14



线性回归：多元线性回归

- 训练数据： $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
 - 其中， $\mathbf{x}_i = [x_{i,1}, x_{i,2}, \dots, x_{i,D}] \in \mathbb{R}^D$ 为代表多元特征的第 i 个数据向量
 - 可以用前 N 天的观看量作为多元特征。
- 学习模型： $f(\mathbf{x}_i) = \sum_{j=1}^D w_j x_{i,j} + w_0 = \mathbf{w}^T \mathbf{x}_i + w_0, (1 \leq i \leq n)$
- 损失函数： $L(\mathbf{w}) = \sum_{i=1}^n \text{loss}(\mathbf{x}_i, y_i) = \sum_{i=1}^n (y_i - (\mathbf{w}^T \mathbf{x}_i + w_0))^2$
- 优化方法：最小二乘法或梯度下降法



北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

$$y = b + wx_1$$

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

$$L' = 0.49k$$

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

1. 如果用前7天的数据作为输入特征，那么特征维度是 [填空1]
2. 如果用前7天的数据作为输入特征，那么该线性模型参数量是 [填空2]
3. 直觉上，你认为前7天（前1天，前2天，...，前7天）的数据中，哪一天的数据对当天观看量的预测最重要？ [填空3]

作答



多元线性模型 ($y = b + \mathbf{w}^T \mathbf{x}$)

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

2017 - 2020 2021
 $L = 0.38k$ $L' = 0.49k$

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 2021
 $L = 0.33k$ $L' = 0.46k$

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

2017 - 2020 2021
 $L = 0.32k$ $L' = 0.46k$

线性模型 *Linear models*



北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

提纲

一、线性回归与梯度下降

二、前馈神经网络

三、卷积神经网络

四、序列数据模型

五、深度学习应用



为何着重讨论线性模型? ($y = b + wx_1$)

- 回顾传统机器学习模型，很多基于线性模型
 - Mean Squared Error (generally for regression)
 - Perceptron
 - Fisher's linear discriminant
 - Support Vector Machines (SVM)
 - Logistic Regression
 - ...
- 前馈神经网络也建立在线性模型基础之上



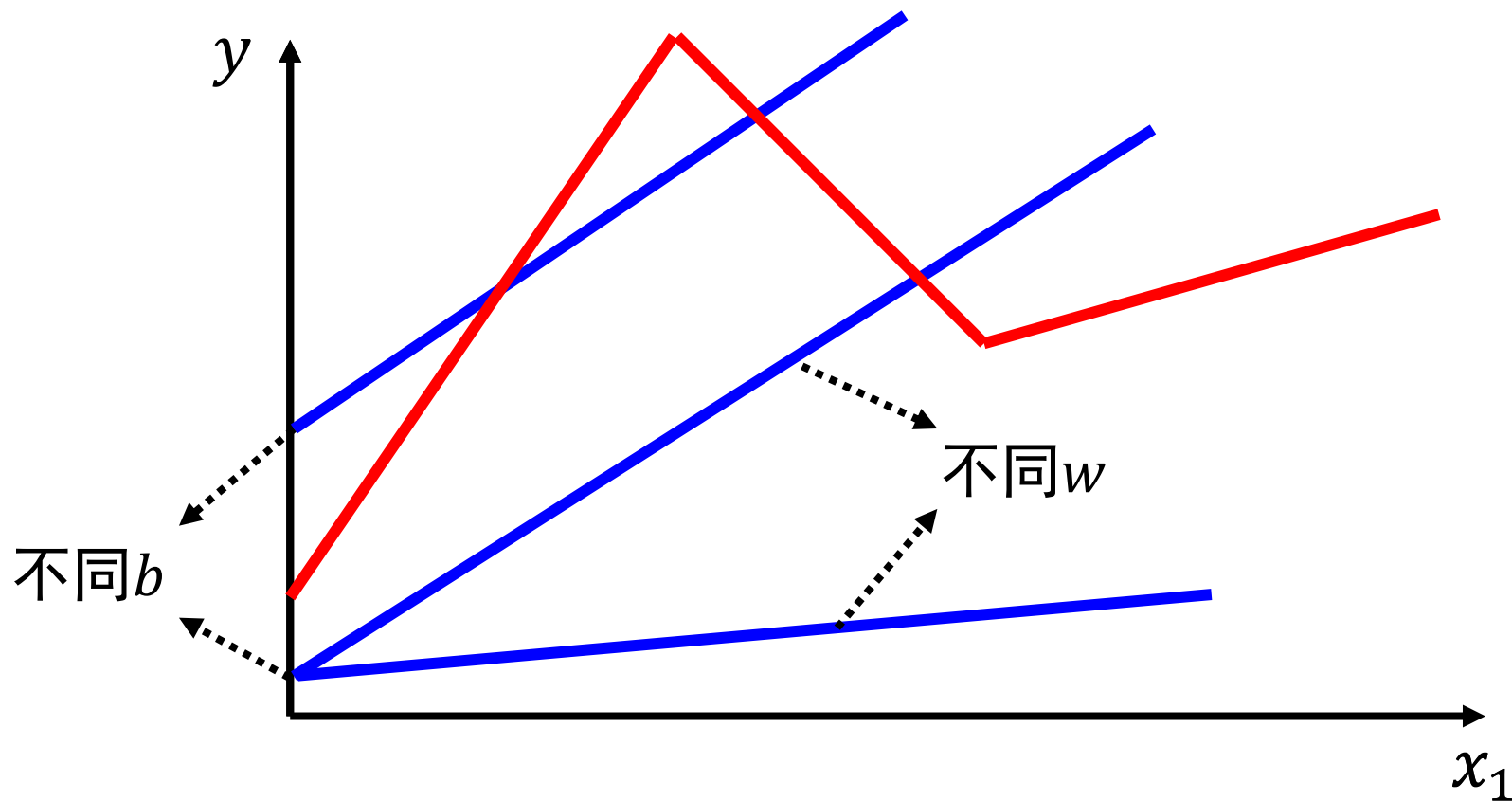
回归分析

- 训练数据: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- 学习模型: $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + w_0, (1 \leq i \leq n)$
 - 扩展线性模型为非线性!
- 损失函数: $L(\mathbf{w}) = \sum_{i=1}^n \text{loss}(\mathbf{x}_i, y_i) = \sum_{i=1}^n (y_i - (\mathbf{w}^T \mathbf{x}_i + w_0))^2$
- 优化方法: 最小二乘法或梯度下降法



北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

线性模型的扩展



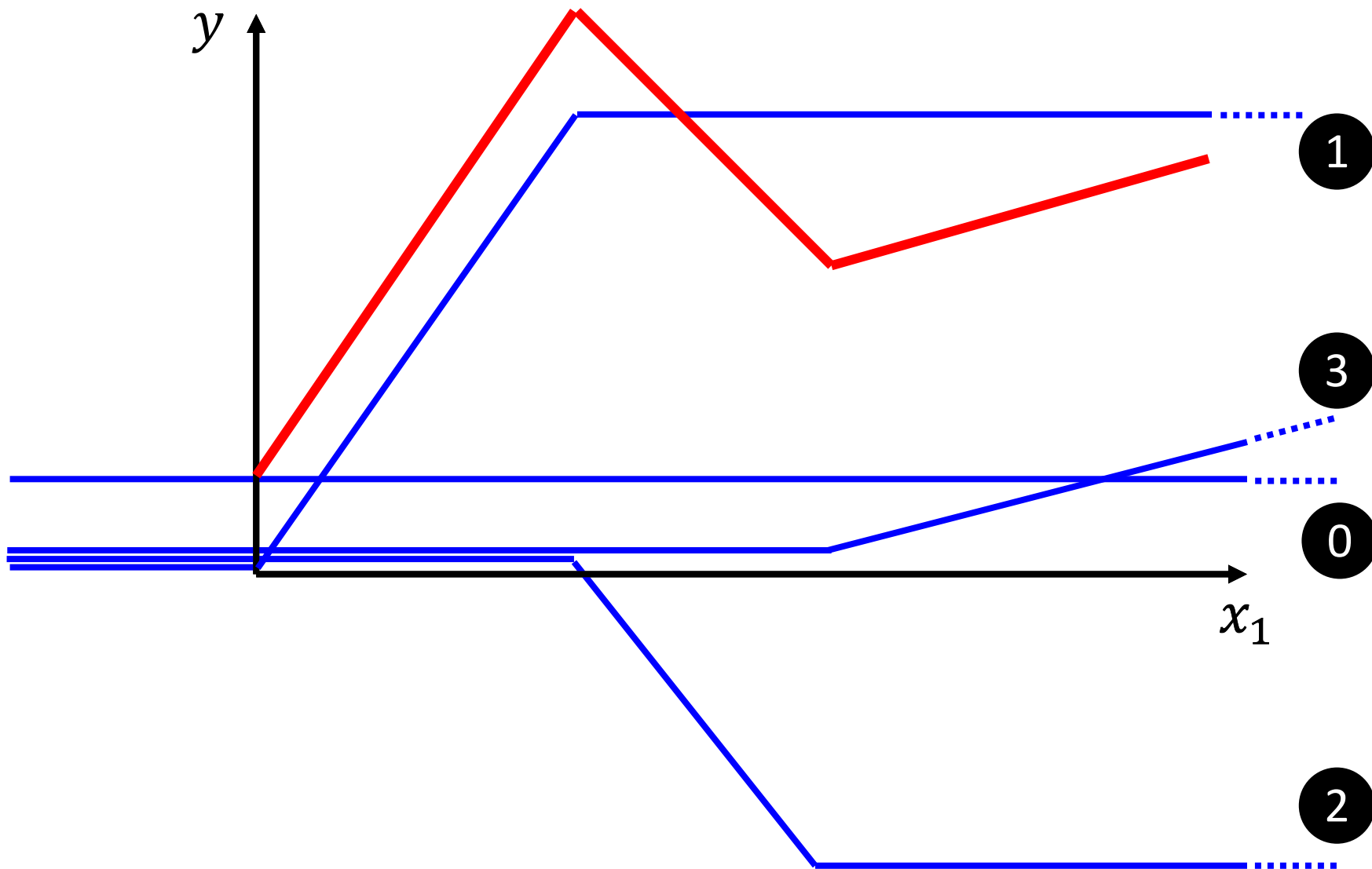
线性模型过于简单(欠拟合)，为了拟合复杂数据需要更加复杂的模型。



北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

线性模型的扩展

red curve = constant + sum of a set of





北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

分段线性函数

= constant + sum of a set of



More pieces require more

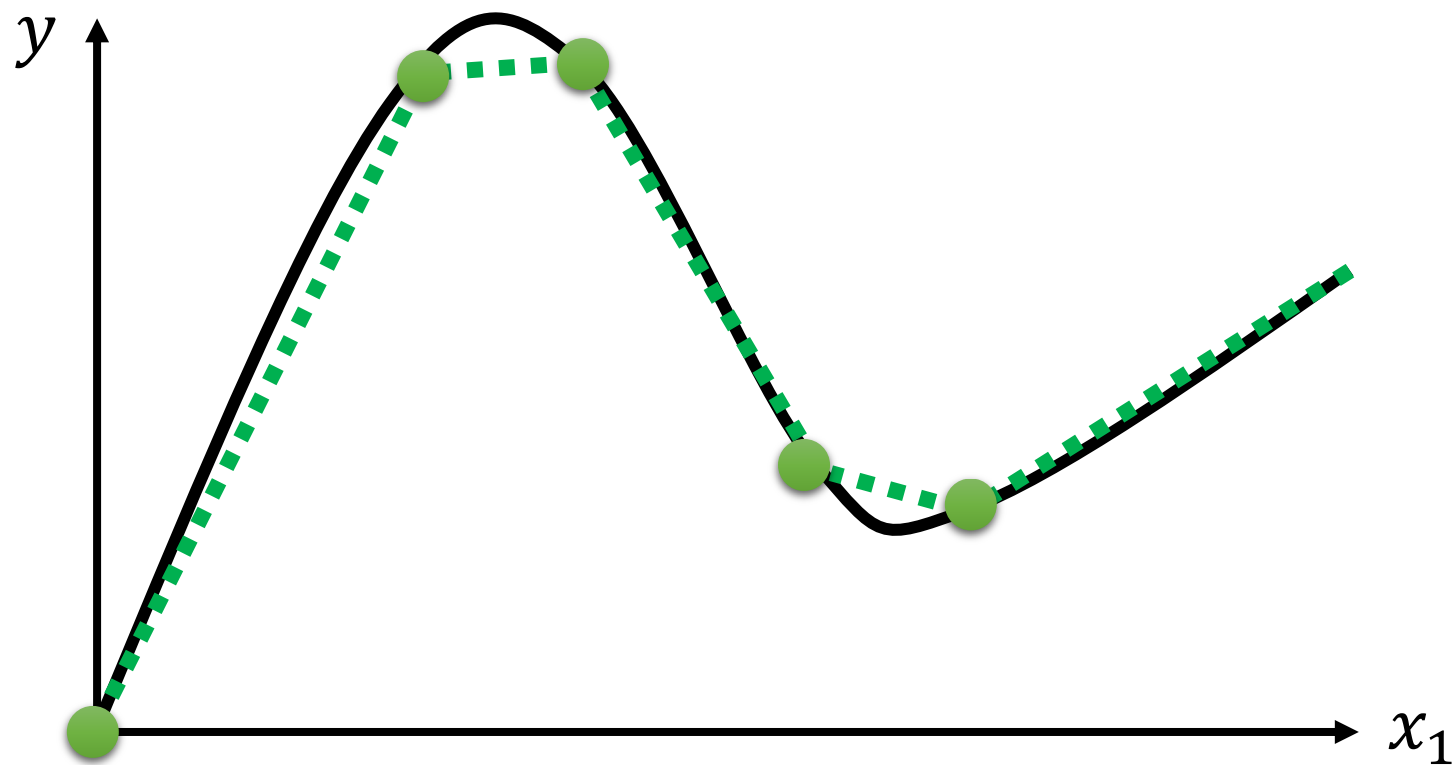




北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

分段线性模型

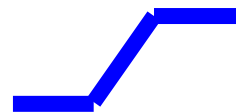
用分段线性模型拟合连续曲线



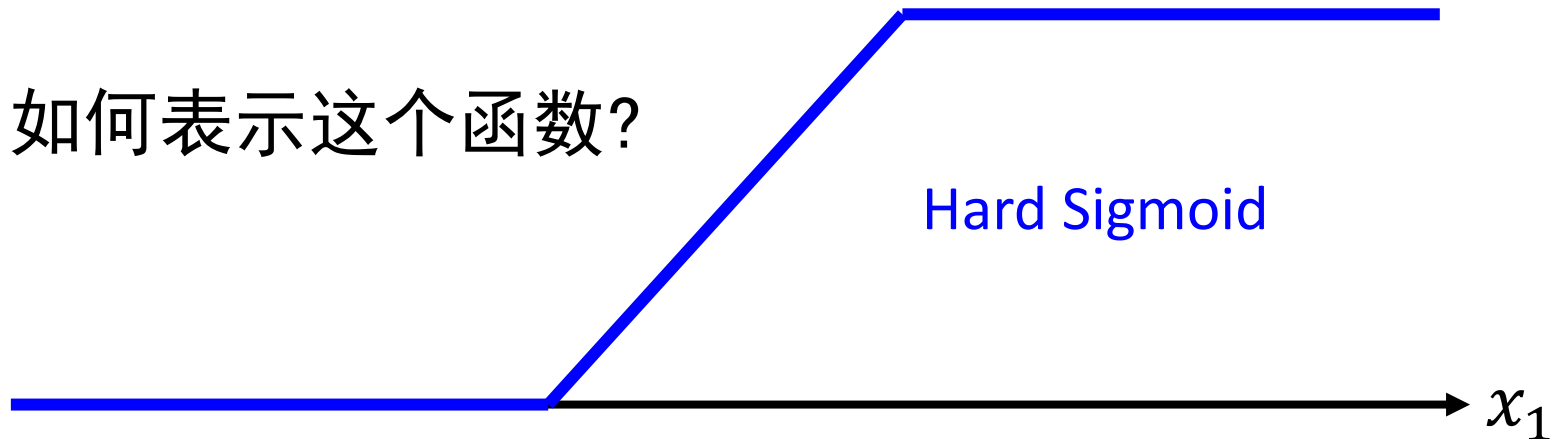
为了更好的近似连续曲线，需要更多分段线性模型



red curve = constant + sum of a set of



如何表示这个函数?

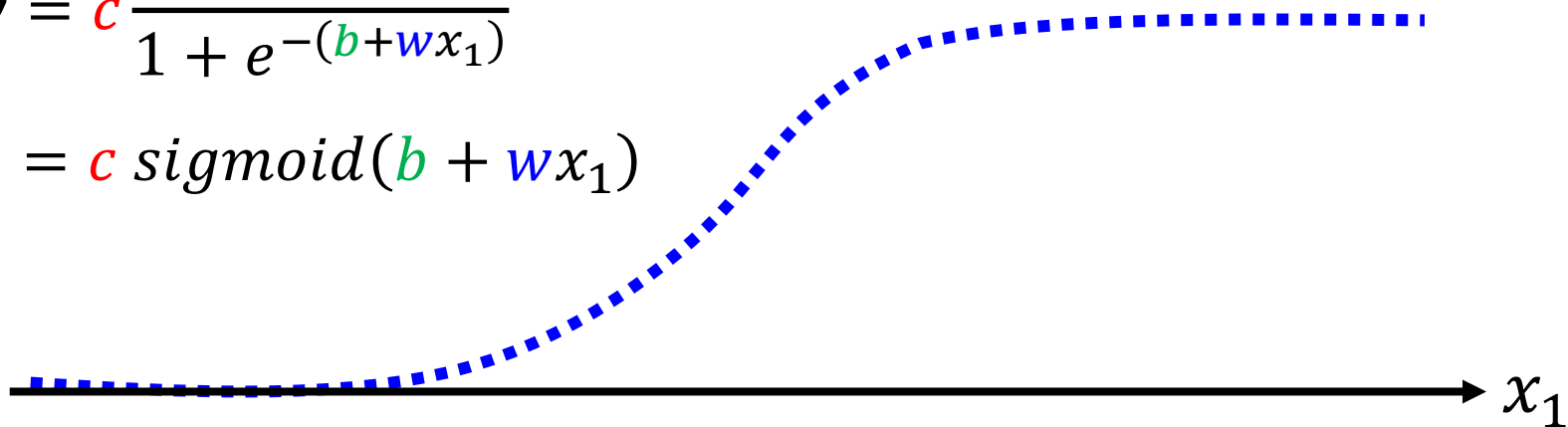


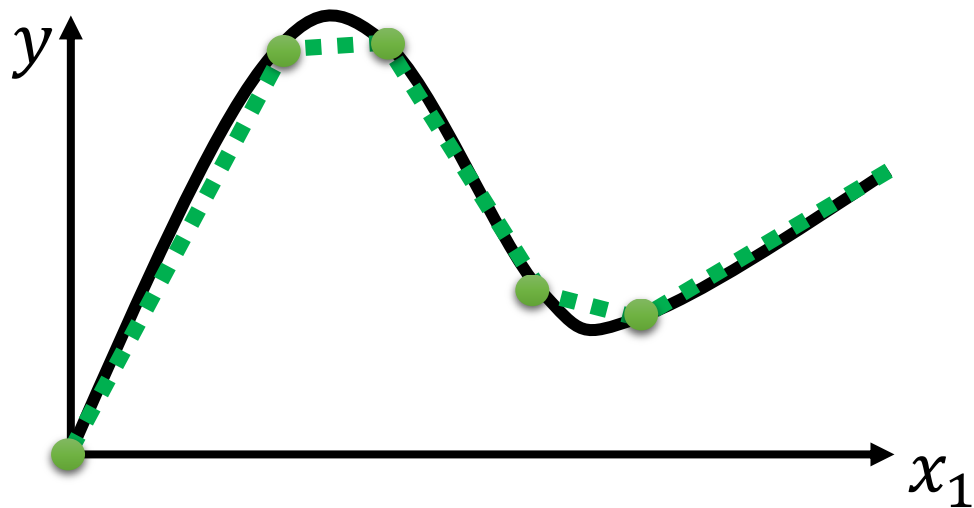
Hard Sigmoid

Sigmoid 函数

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

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





$$\text{分段线性函数} = \text{常量} + \sum \text{ / }$$

如果激活函数是如下简单线性模型（而非分段线性），是否还能近似复杂非线性函数 [填空1]？为什么 [填空2]？

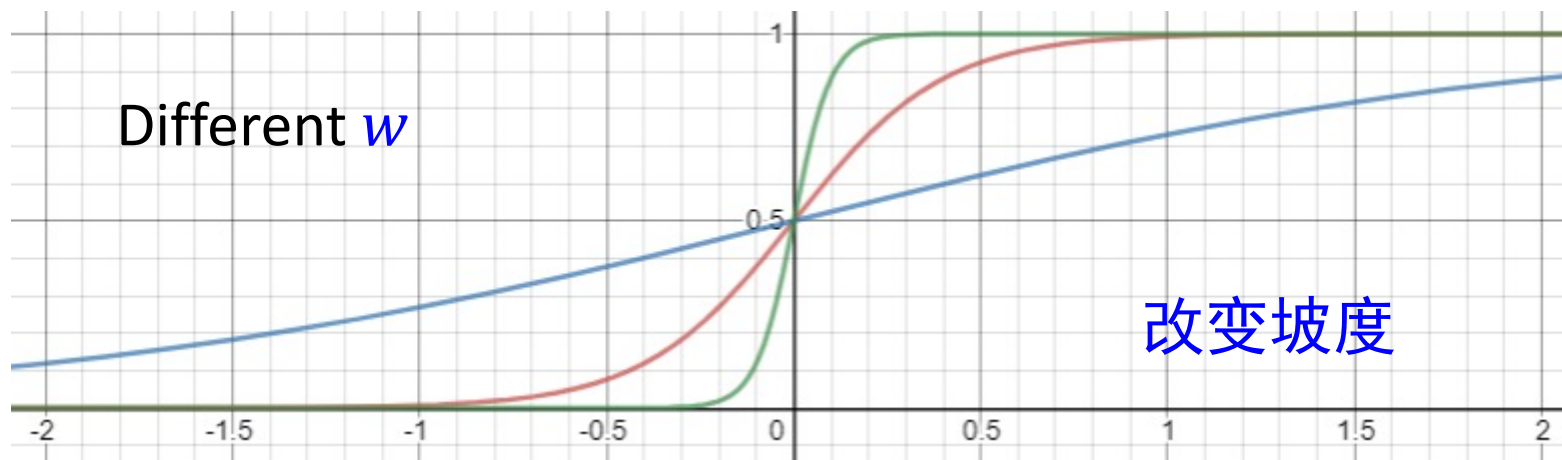
$$\text{常量} + \sum \text{ / }$$

作答

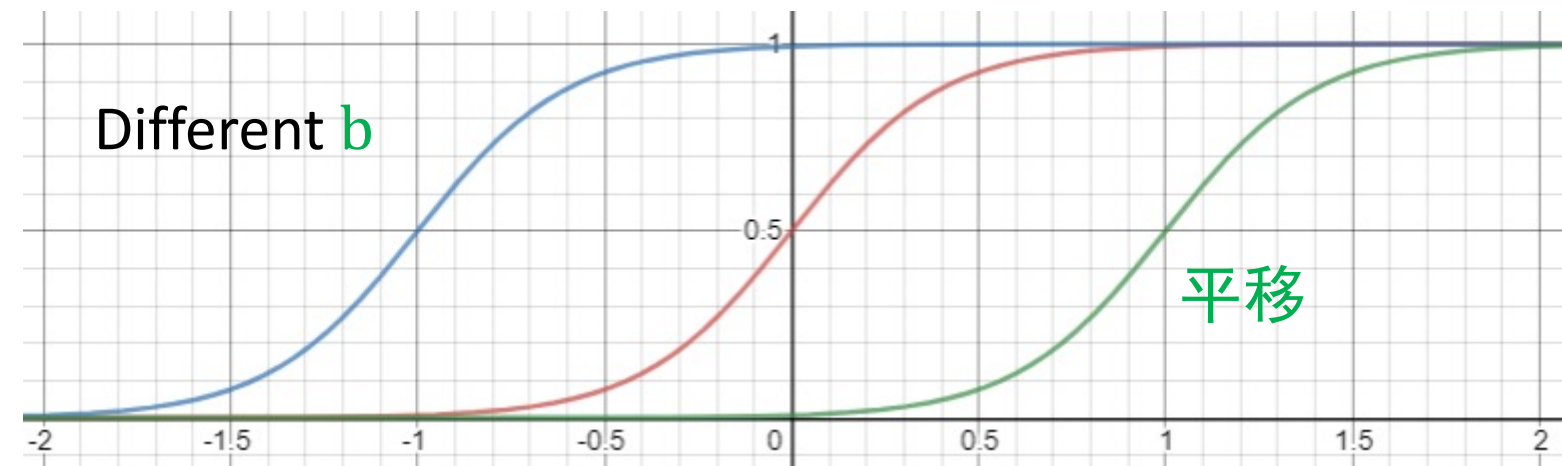


北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

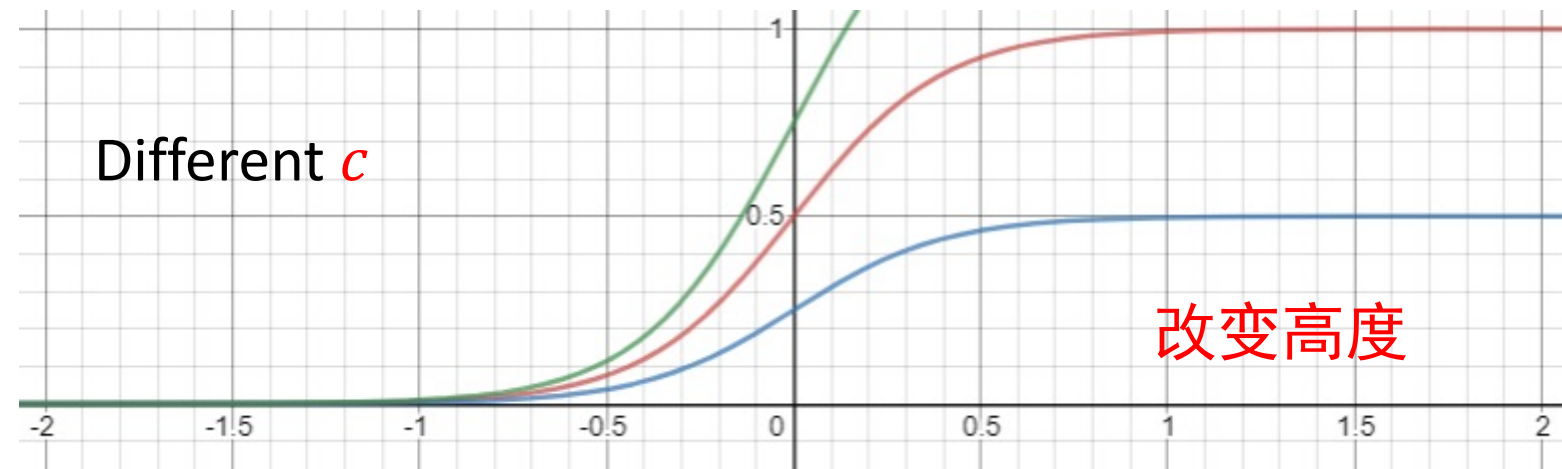
Different w



Different b

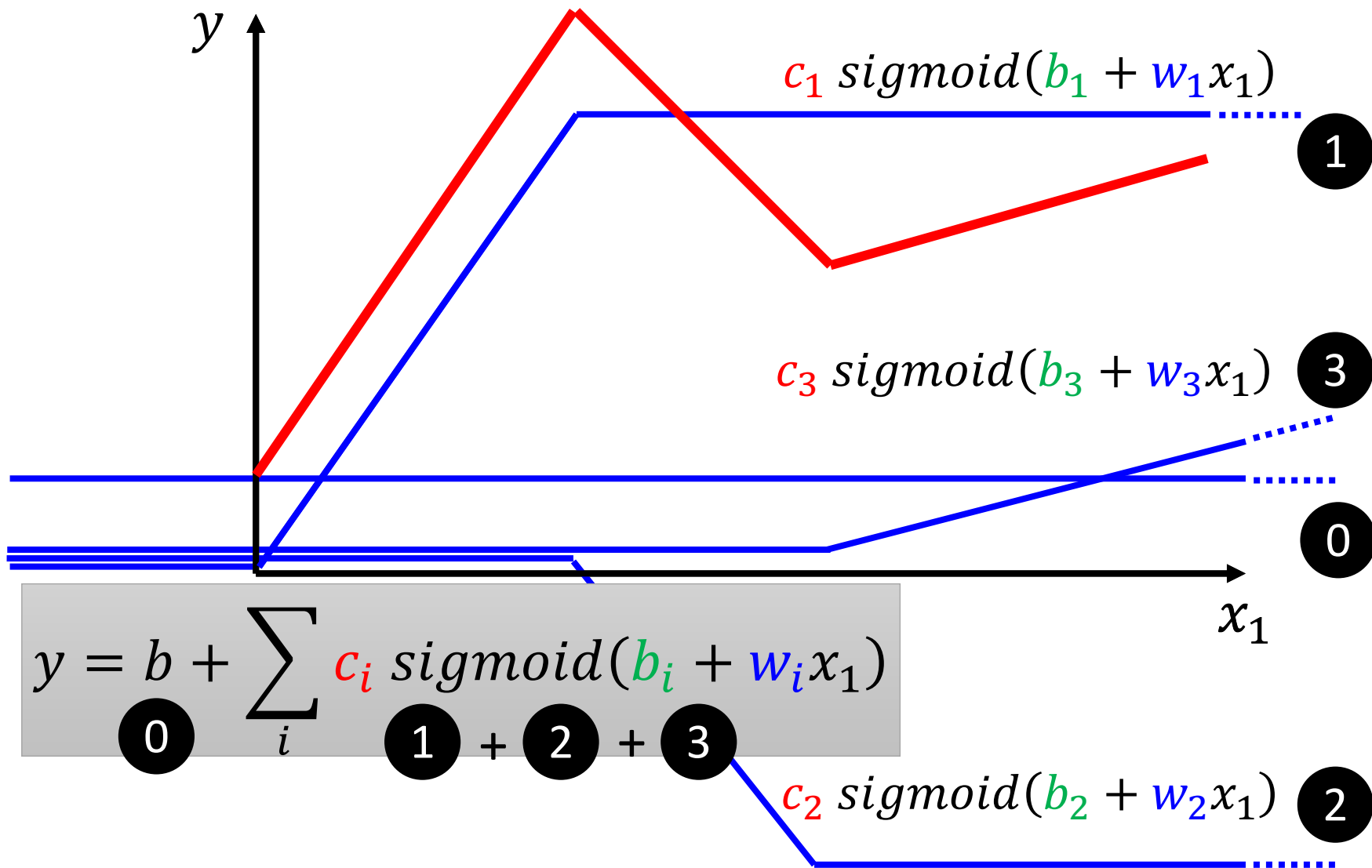


Different c





red curve = sum of a set of  + constant





新模型：更多模型组合

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

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

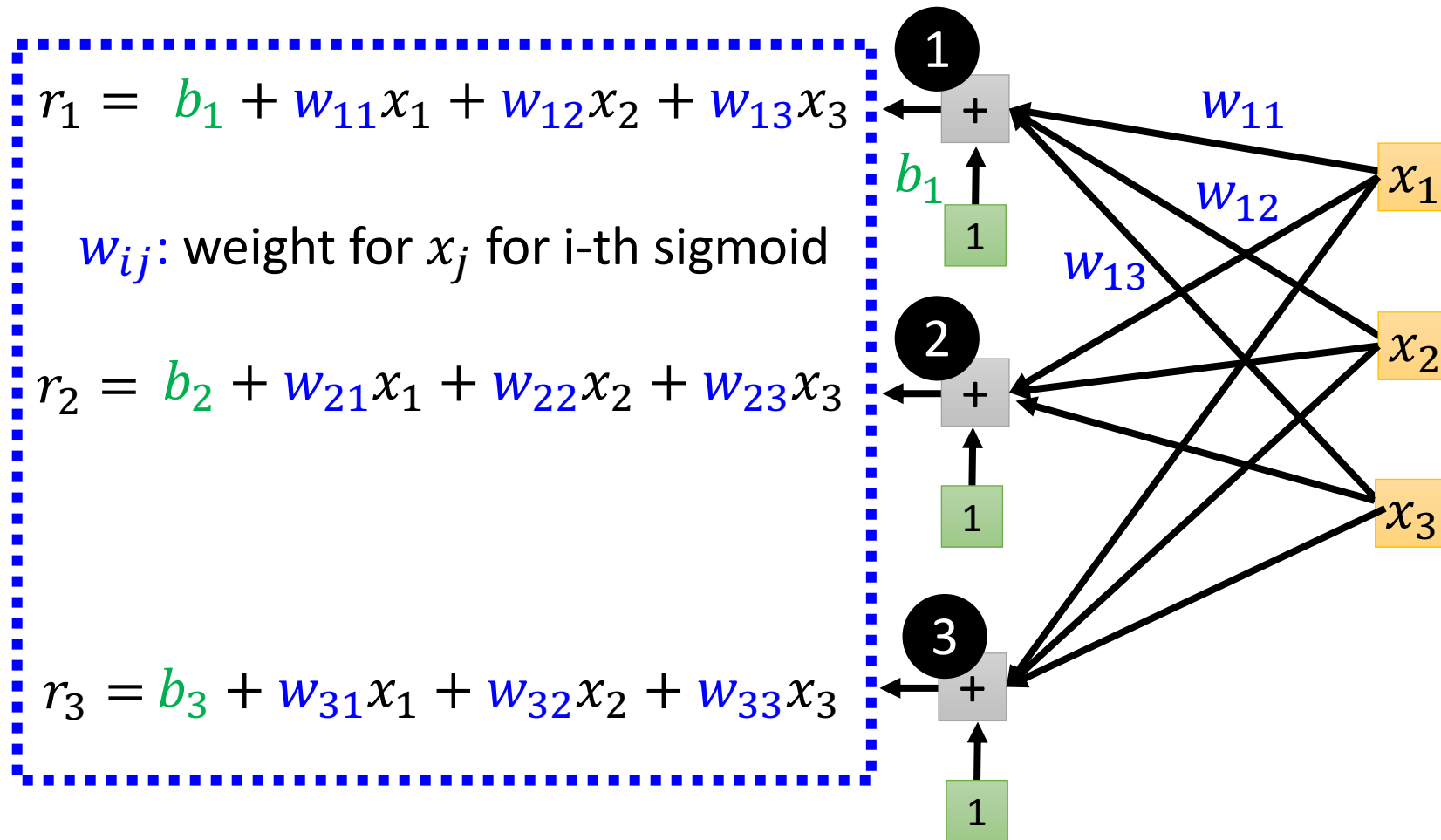
$$y = b + \underline{\sum_j w_j x_j}$$

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



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

$j: 1, 2, 3$ 特征个数
 $i: 1, 2, 3$ sigmoid个数





$$y = b + \sum_i c_i \operatorname{sigmoid} \left(b_i + \sum_j w_{ij} x_j \right) \quad \begin{array}{l} i: 1, 2, 3 \\ j: 1, 2, 3 \end{array}$$

$$r_1 = b_1 + w_{11}x_1 + w_{12}x_2 + w_{13}x_3$$

$$r_2 = b_2 + w_{21}x_1 + w_{22}x_2 + w_{23}x_3$$

$$r_3 = b_3 + w_{31}x_1 + w_{32}x_2 + w_{33}x_3$$

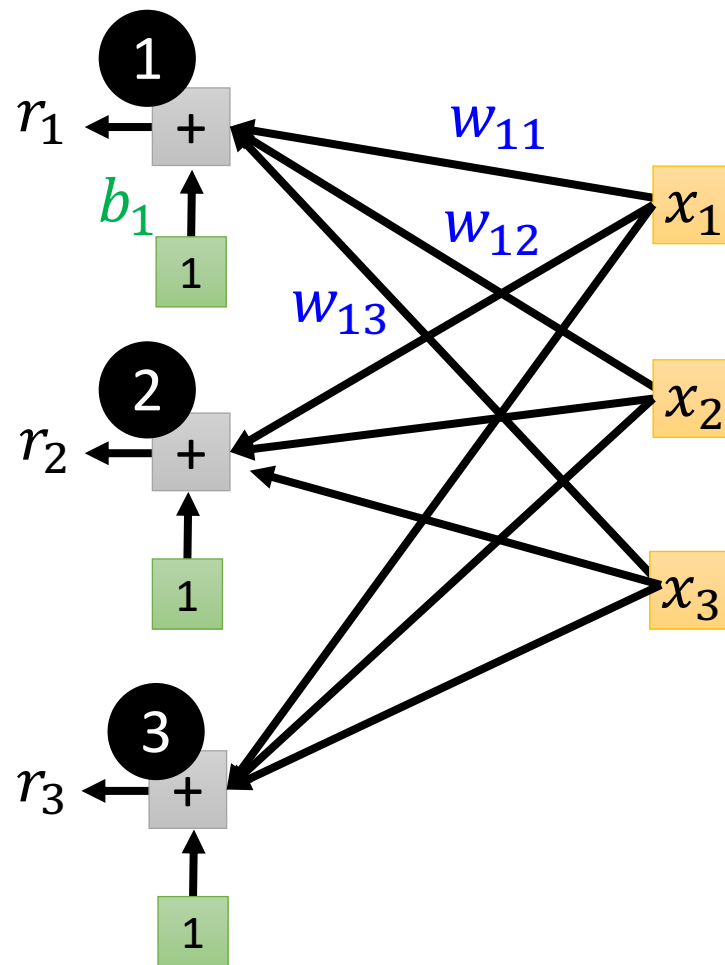
$$\begin{bmatrix} r_1 \\ r_2 \\ r_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} + \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

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



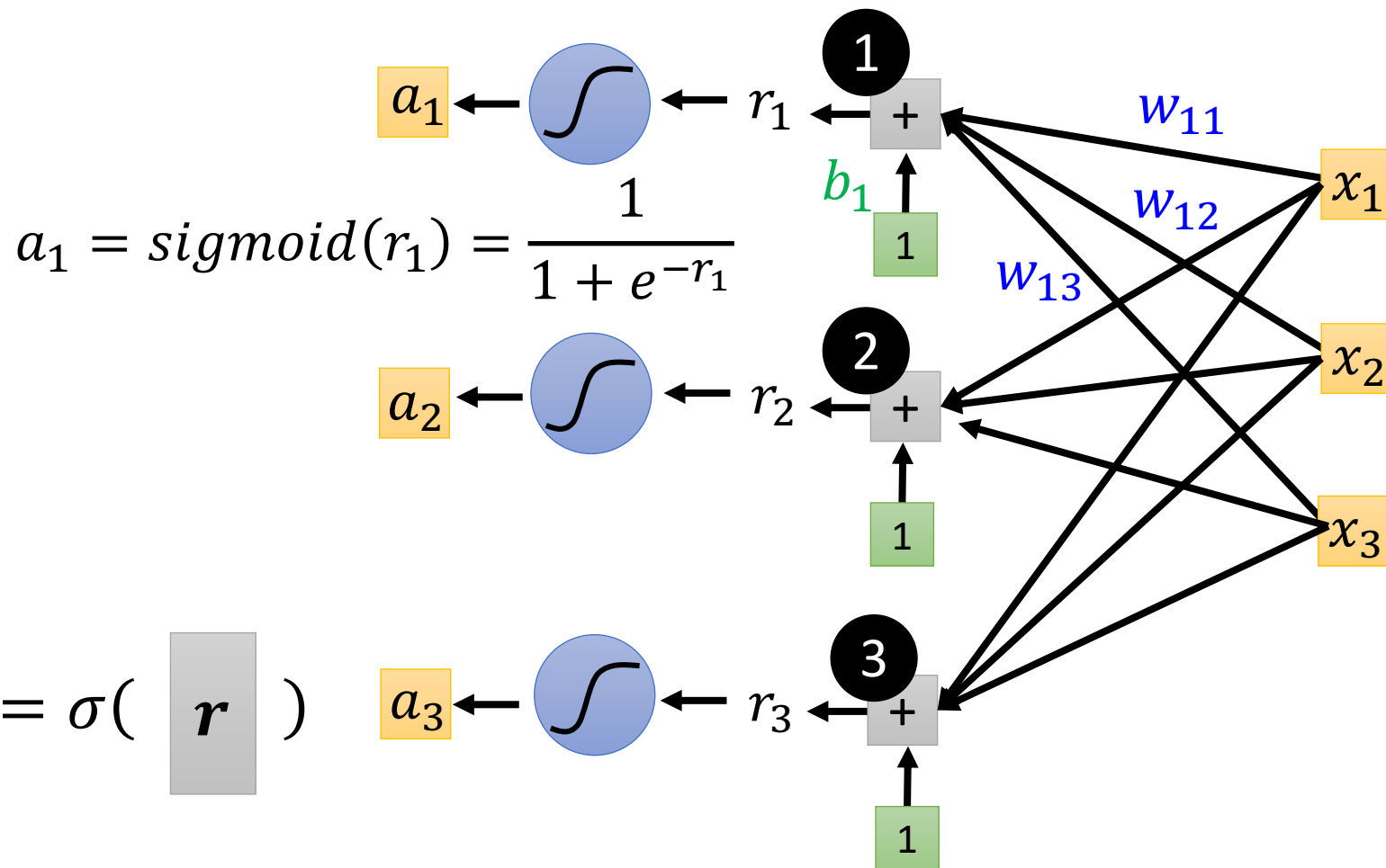
$$y = b + \sum_i c_i \operatorname{sigmoid} \left(b_i + \sum_j w_{ij} x_j \right) \quad \begin{array}{l} i: 1, 2, 3 \\ j: 1, 2, 3 \end{array}$$

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



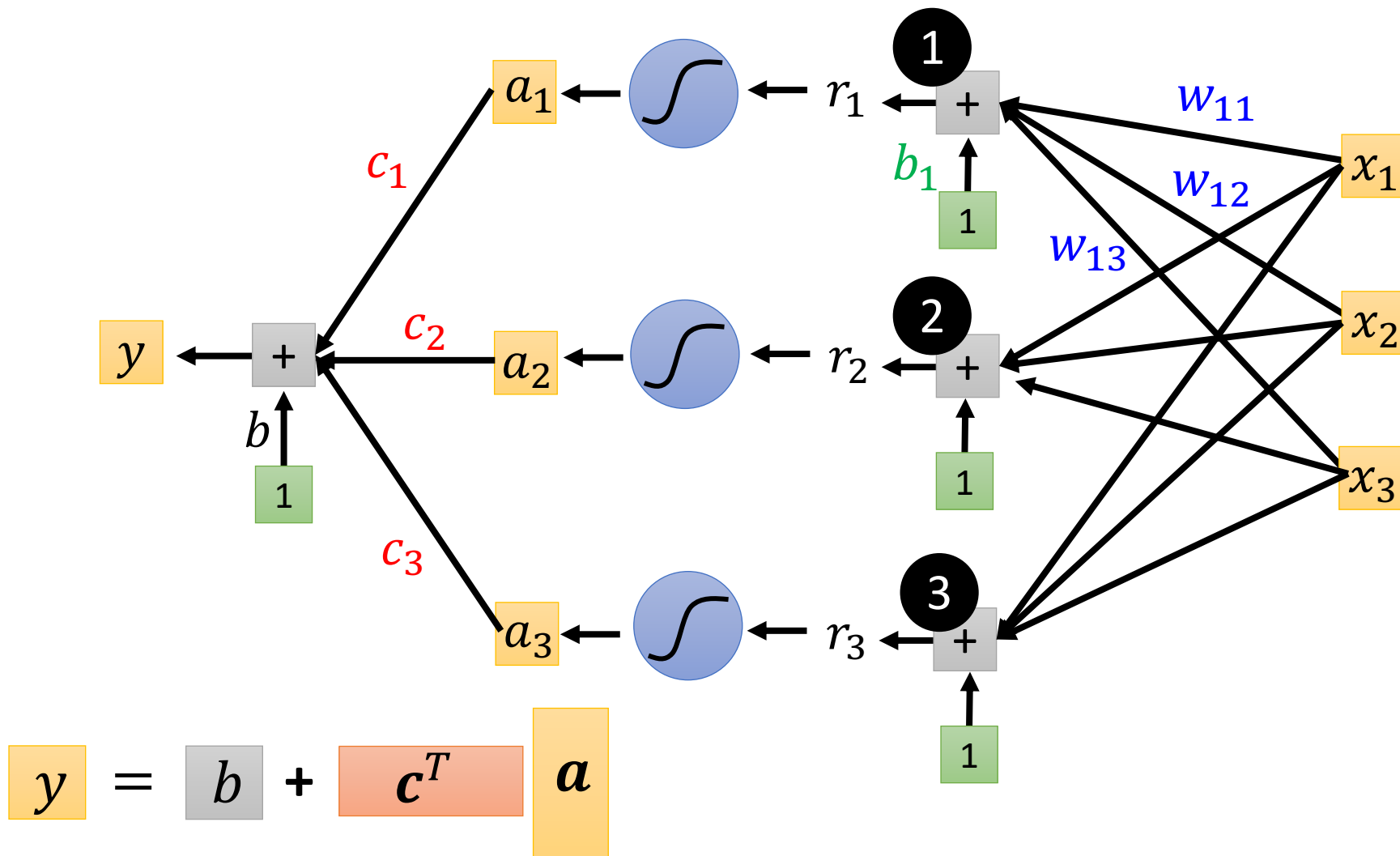


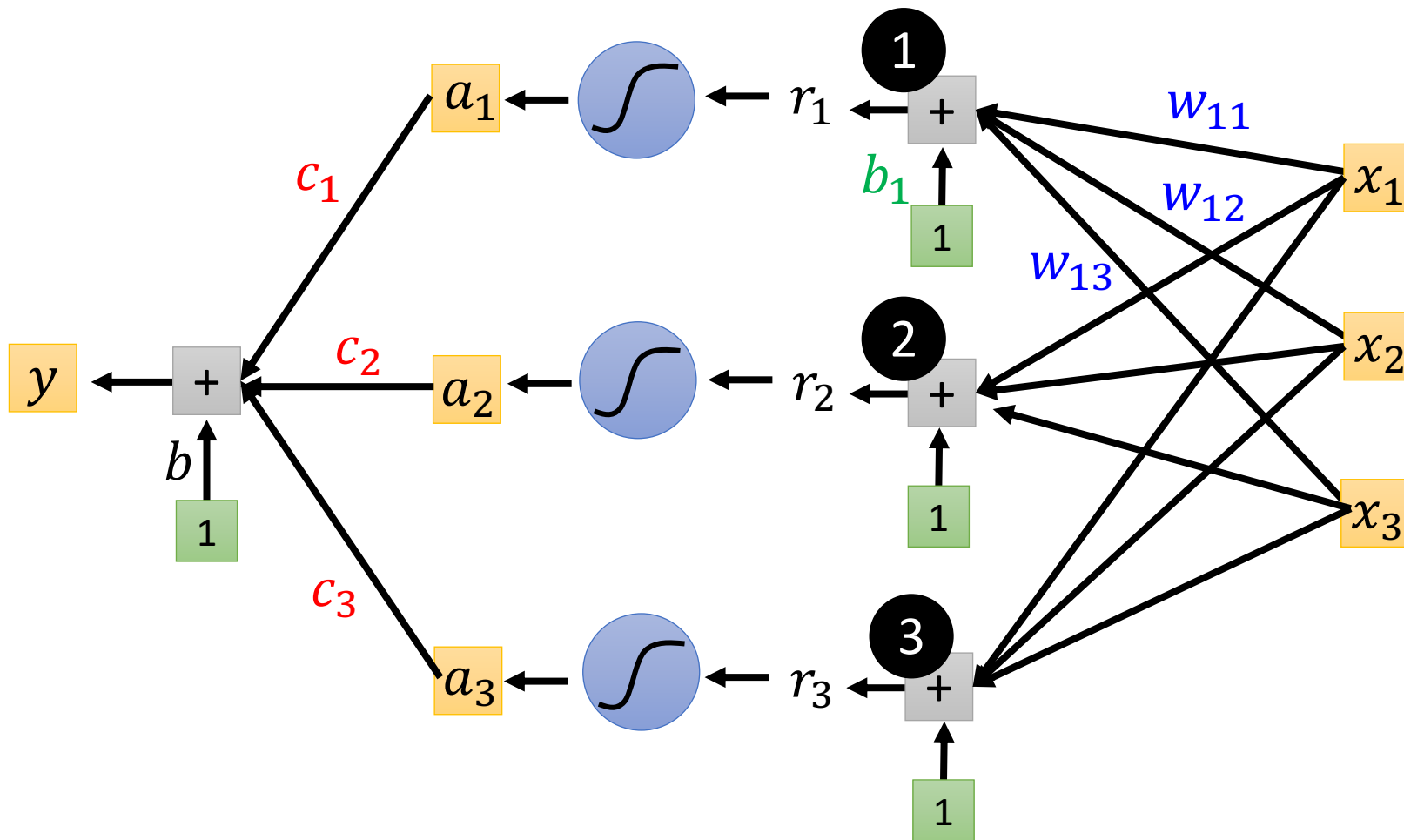
$$y = b + \sum_i c_i \text{sigmoid} \left(b_i + \sum_j w_{ij} x_j \right) \quad \begin{matrix} i: 1,2,3 \\ j: 1,2,3 \end{matrix}$$



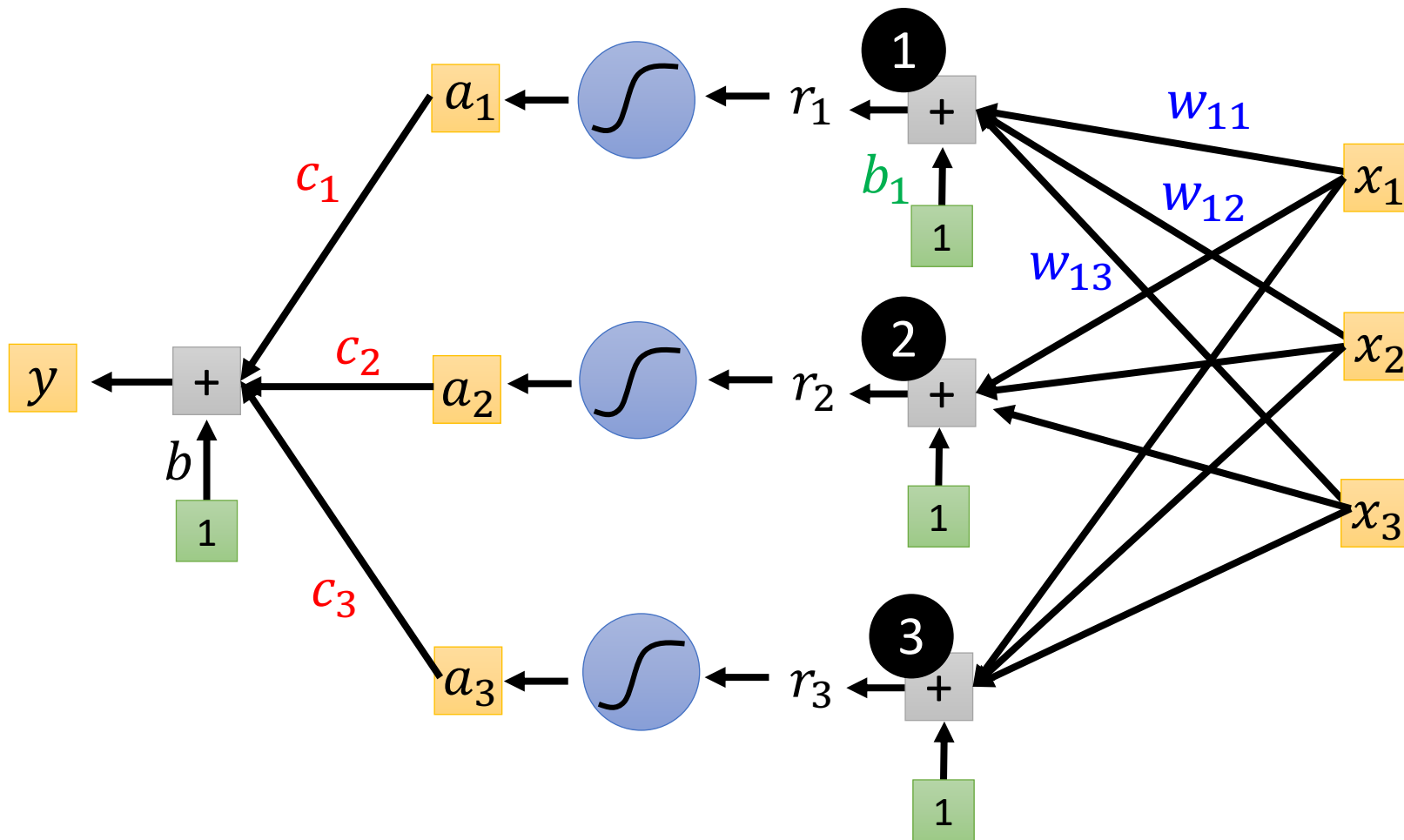


$$y = b + \sum_i c_i \operatorname{sigmoid} \left(b_i + \sum_j w_{ij} x_j \right) \quad \begin{array}{l} i: 1, 2, 3 \\ j: 1, 2, 3 \end{array}$$





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

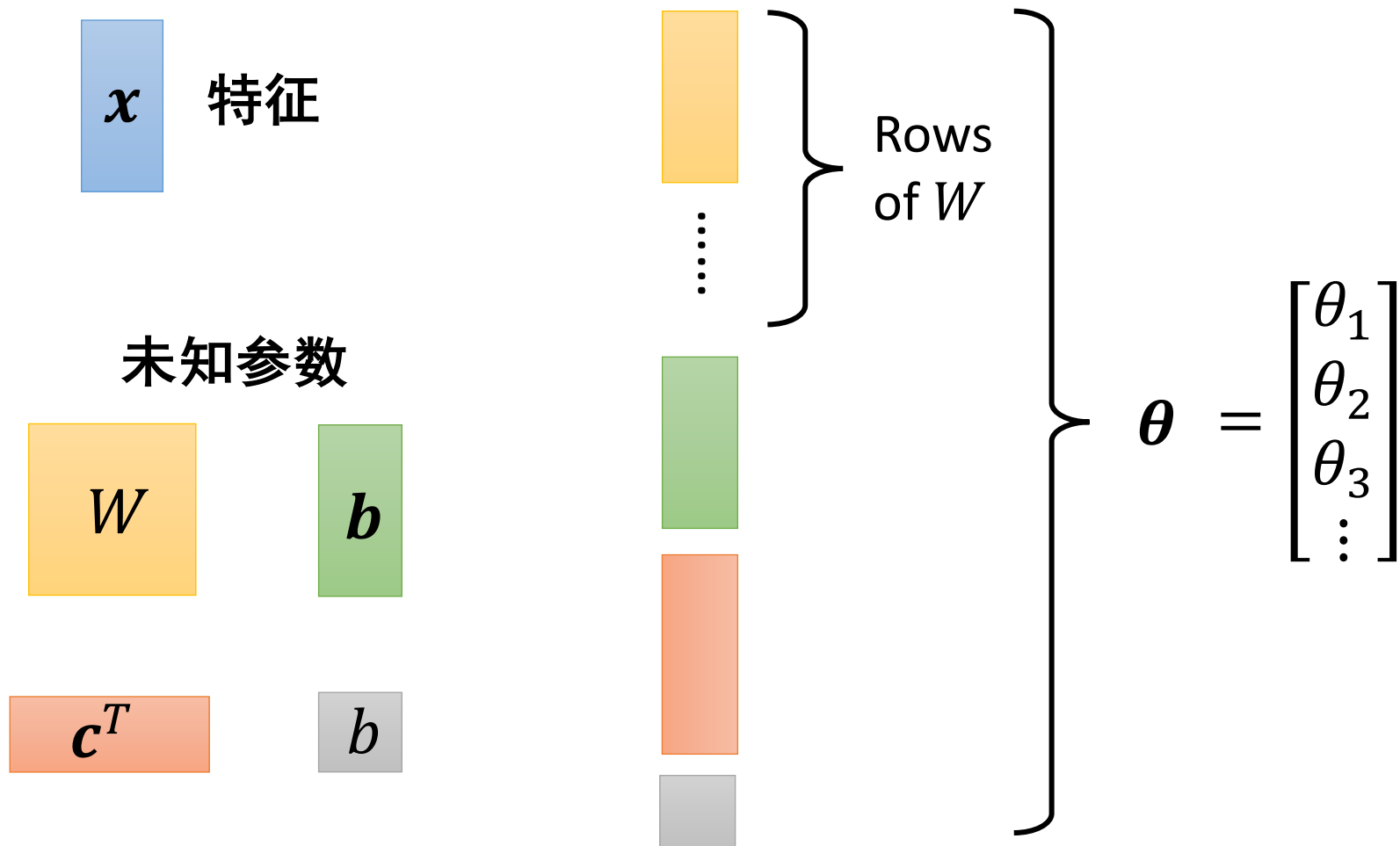


$$y = b + c^T \sigma(b + Wx)$$



带有未知参数的函数

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





回归分析：非线性回归

- 训练数据: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- 学习模型: $f(\mathbf{x}_i) = b + \mathbf{c}^T \sigma(\mathbf{b} + W\mathbf{x}_i), (1 \leq i \leq n)$
- 损失函数: $L(\boldsymbol{\theta}) = \sum_{i=1}^n \text{loss}(\mathbf{x}_i, y_i) = \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$
- 优化方法: 梯度下降法



优化方法：梯度下降

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

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

➤ (随机) 初始化参数 θ^0

$$\underset{\text{梯度}}{g} = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} |_{\theta=\theta^0} \\ \frac{\partial L}{\partial \theta_2} |_{\theta=\theta^0} \\ \vdots \end{bmatrix}$$

$$g = \nabla L(\theta^0)$$

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

$$\theta^1 \leftarrow \theta^0 - \eta g$$



优化方法：梯度下降

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

➤ (随机) 初始化参数 θ^0

➤ 计算梯度 $g = \nabla L(\theta^0)$

$$\theta^1 \leftarrow \theta^0 - \eta g$$

➤ 计算梯度 $g = \nabla L(\theta^1)$

$$\theta^2 \leftarrow \theta^1 - \eta g$$

➤ 计算梯度 $g = \nabla L(\theta^2)$

$$\theta^3 \leftarrow \theta^2 - \eta g$$



回归分析

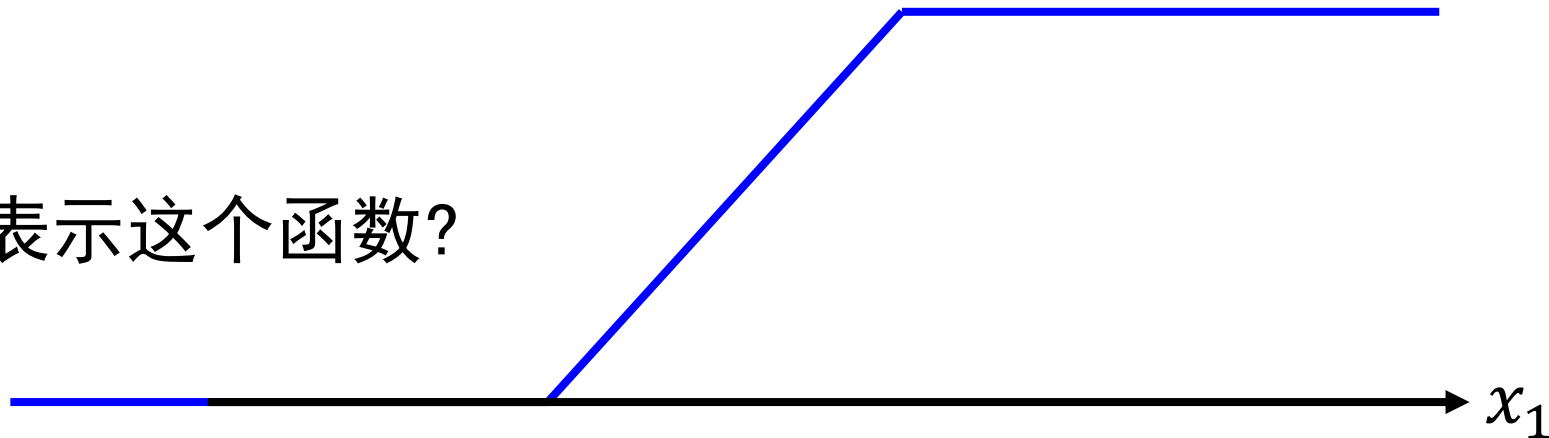
- 训练数据: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- 学习模型: $f(\mathbf{x}_i) = b + \mathbf{c}^T \sigma(\mathbf{b} + W\mathbf{x}_i), (1 \leq i \leq n)$
 - 通过改变激活函数, 改进学习模型!
- 损失函数: $L(\boldsymbol{\theta}) = \sum_{i=1}^n \text{loss}(\mathbf{x}_i, y_i) = \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$
- 优化方法: 梯度下降法



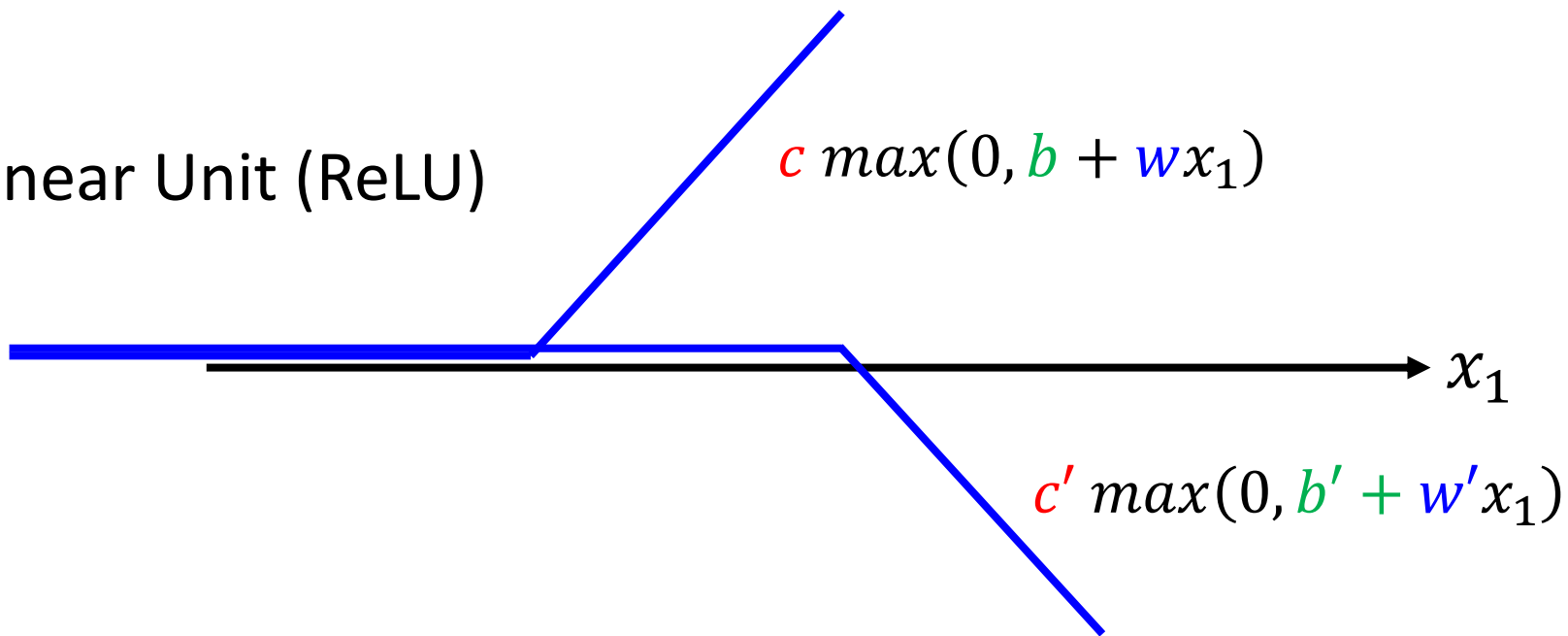
北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

Sigmoid \rightarrow ReLU

如何表示这个函数？



Rectified Linear Unit (ReLU)





Sigmoid \rightarrow ReLU

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

激活函数 Activation function

$$y = b + \sum_{2i} c_i \max \left(0, b_i + \sum_j w_{ij} x_j \right)$$

如何选择激活函数？



北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

实验结果

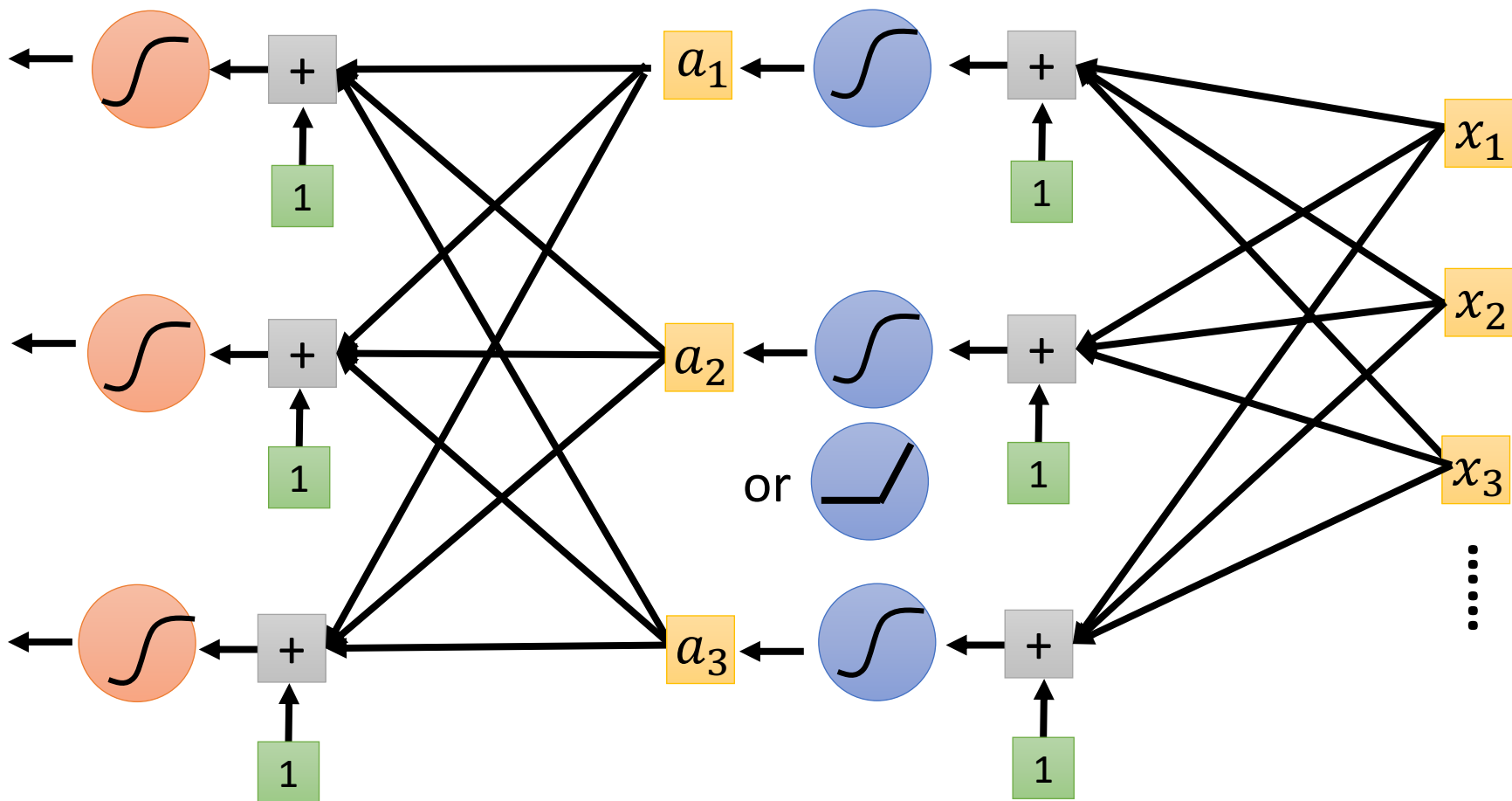
$$y = b + \sum_{\textcolor{red}{i}} \textcolor{red}{c}_i \max \left(0, \textcolor{green}{b}_i + \sum_j \textcolor{blue}{w}_{ij} x_j \right)$$

	linear
2017 – 2020	0.32k
2021	0.46k



回归分析

- 训练数据: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- 学习模型: $f(\mathbf{x}_i) = b + \mathbf{c}^T \sigma(\mathbf{b} + W\mathbf{x}_i), (1 \leq i \leq n)$
 - 通过增加模型层数，继续改进学习模型！
- 损失函数: $L(\boldsymbol{\theta}) = \sum_{i=1}^n \text{loss}(\mathbf{x}_i, y_i) = \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$
- 优化方法: 梯度下降法



$$\mathbf{a}' = \sigma(\mathbf{b}' + \mathbf{W}' \mathbf{a}) \quad \mathbf{a} = \sigma(\mathbf{b} + \mathbf{W} \mathbf{x})$$



回归分析

- 训练数据: $\mathcal{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$
- 学习模型: $f(\mathbf{x}_i) = b + \mathbf{c}^T \sigma(\mathbf{b}' + W'(\sigma(\mathbf{b} + W\mathbf{x}_i))), (1 \leq i \leq n)$
- 损失函数: $L(\boldsymbol{\theta}) = \sum_{i=1}^n \text{loss}(\mathbf{x}_i, y_i) = \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$
- 优化方法: 梯度下降法



实验结果

- 多个隐含层模型的损失
 - 每层100个 ReLU 神经元
 - 输入特征是过去56天的观看量

	1 layer
2017 – 2020	0.28k
2021	0.43k



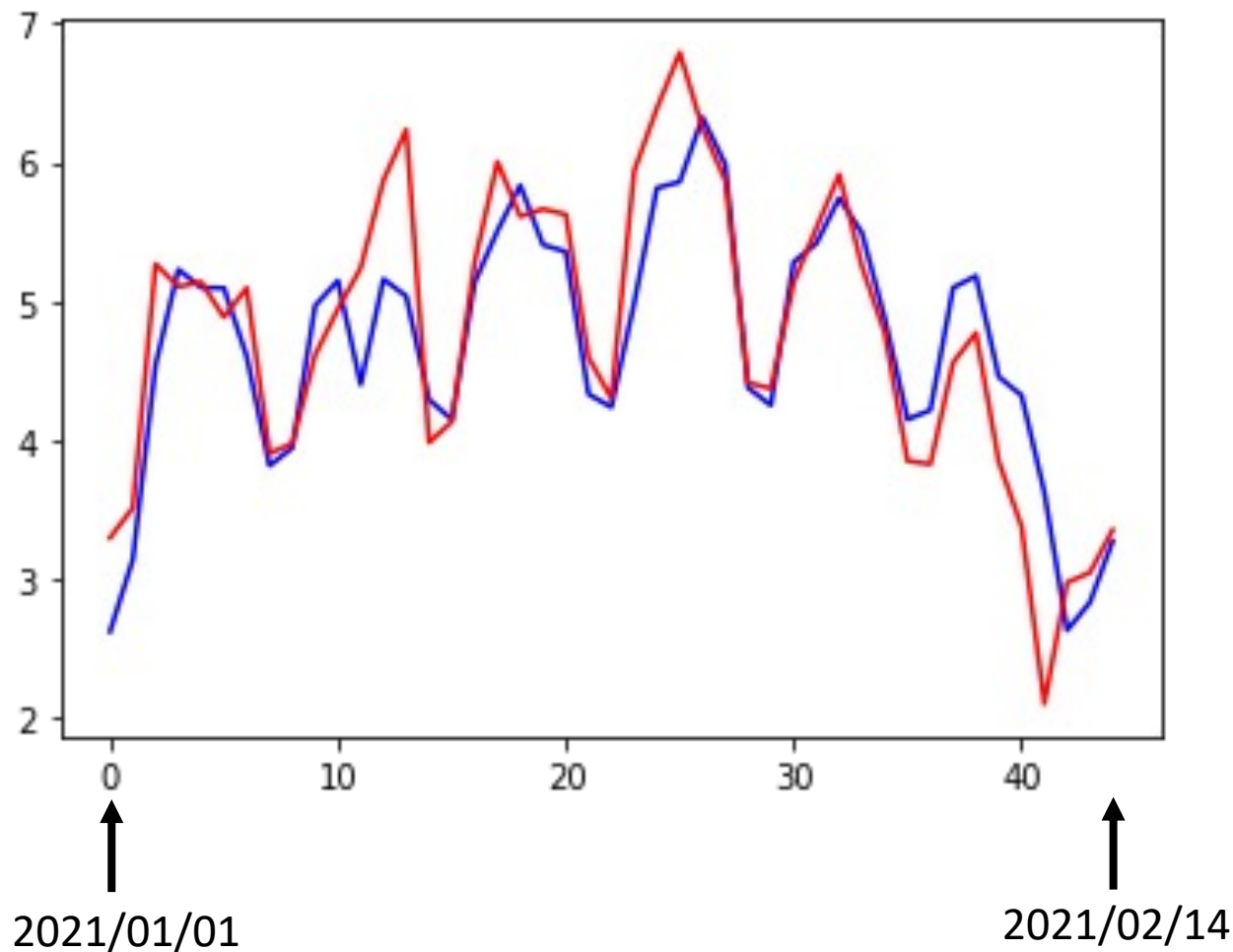
北京航空航天大学
COLLEGE OF SOFTWARE
BEIHANG UNIVERSITY 软件学院

3 layers

Red: 真实观看量

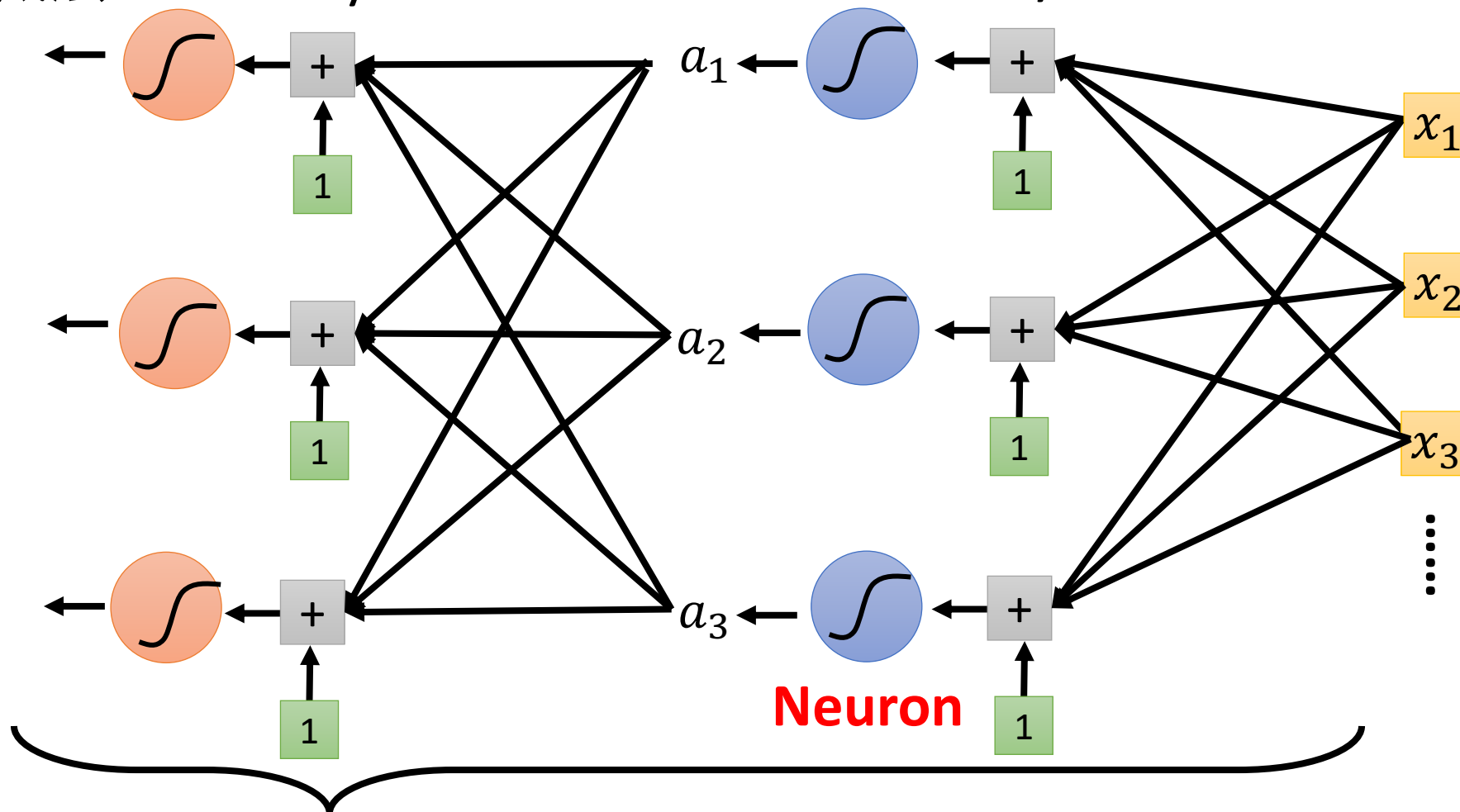
blue: 模型预测观看量

观看量
(k)



隐藏层hidden layer

隐藏层hidden layer



神经网络 Neural Network

Many layers means **Deep** ➡ 深度学习 Deep Learning



前馈神经网络

- 输入层、输出层和至少一层的隐藏层构成。网络中各个隐藏层中神经元可接收相邻前序隐藏层中所有神经元传递而来的信息，经过加工处理后将信息输出给相邻后续隐藏层中所有神经元。
- 各个神经元接受前一级的输入，并输出到下一级，模型中没有反馈
- 层与层之间通过“全连接”进行链接，即两个相邻层之间的神经元完全成对连接，但层内的神经元不相互连接。
- 也被称为全连接网络，或多层感知机。