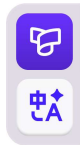
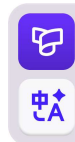
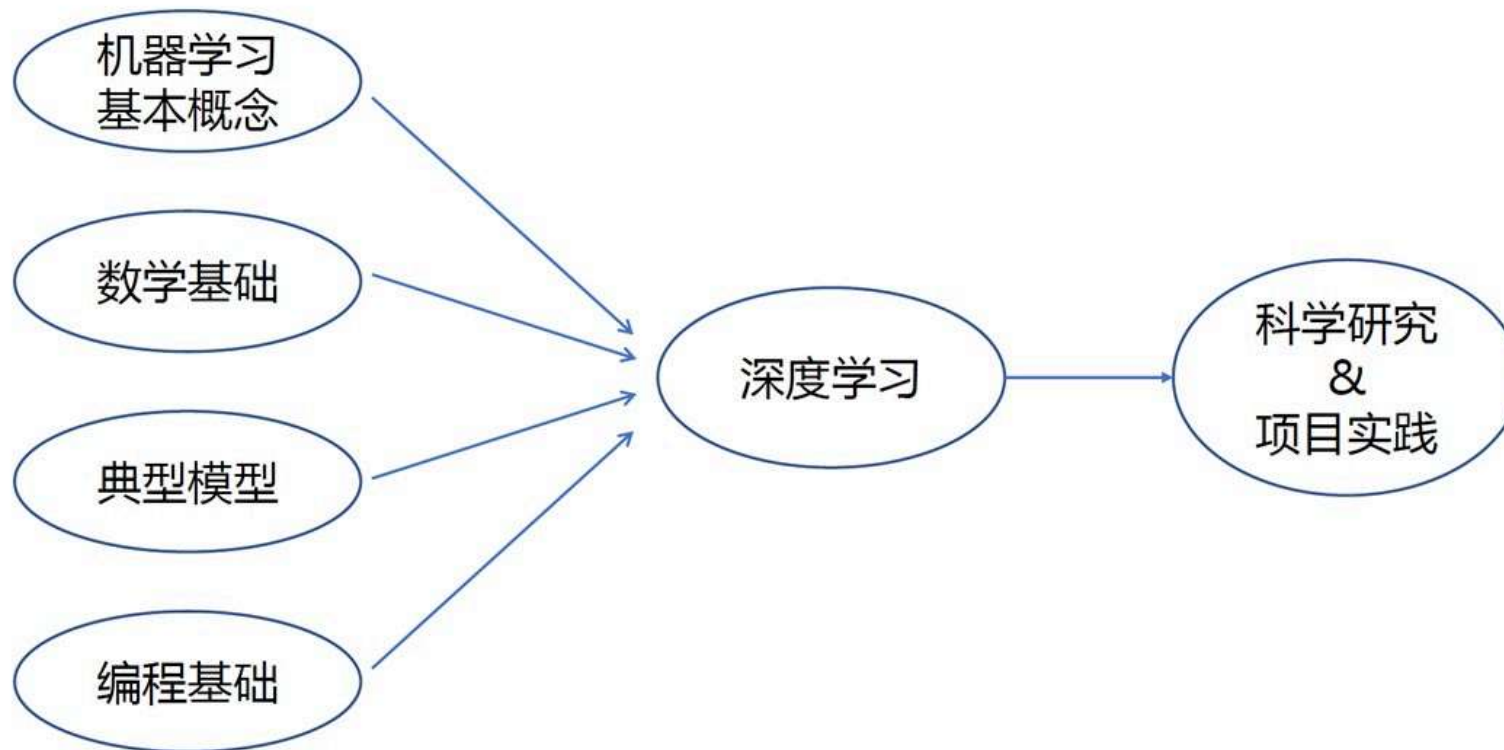


第二讲 基础知识



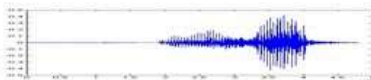
深度学习基础 & 数据基础



深度学习基础

➤ 深度(机器)学习 \approx 构建一个映射函数


◆ 语音识别

$f(\text{  }) = \text{“你好”}$

◆ 图像识别

$f(\text{  }) = \text{“猫”}$

◆ 围棋

$f(\text{  }) = \text{“5-5” (落子位置)}$

◆ 对话系统

$f(\text{“你好”}) = \text{“今天天气真不错”}$

用户输入

机器回应

[Neural networks and deep learning](#)

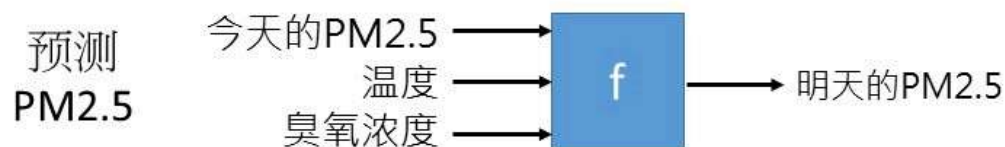
<https://udlbook.github.io/udlbook/>



深度学习的类别

➤ 深度学习的类别(Different types of Functions)

◆ 回归(Regression): 函数的输出是标量.



◆ 分类(Classification): 给定选项(类), 函数输出正确的选项



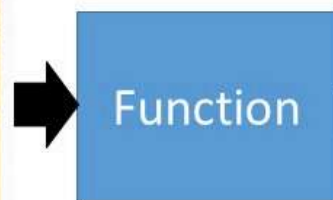
深度学习的类别

➤ 深度学习的类别(Different types of Functions)

◆ **分类(Classification)**: 给定选项(类), 函数输出正确的选项



下棋



棋盘上的一个位置

下一步落子位置

每个位置是一个分类
(19 x 19 类)

从观看视频人数开始

我们要找的函数…

$y = f(\text{2月27日的观看人数})$

影片	流量来源	地理位置	观众年龄	观众性别	日期	订阅状态	订阅来源	播放清单
日期 ↓					喜欢的人数	订阅人数	观看次数	
2021年1月26日					54 4.9%	69 5.5%	6,788	5.2%
2021年1月27日					60 5.4%	71 5.6%	6,242	4.7%
2021年1月28日					36 3.2%	63 5.0%	5,868	4.5%
2021年1月29日					27 2.4%	40 3.2%	4,413	3.4%
2021年1月30日					40 3.6%	40 3.2%	4,372	3.3%
2021年1月31日					47 4.2%	51 4.0%	5,135	3.9%
2021年2月1日					61 5.5%	29 2.3%	5,527	4.2%
2021年2月2日					49 4.4%	43 3.4%	5,911	4.5%
2021年2月3日					26 2.3%	44 3.5%	5,248	4.0%
2021年2月4日					43 3.9%	33 2.6%	4,771	3.6%
2021年2月5日					45 4.0%	49 3.9%	3,850	2.9%
2021年2月6日					29 2.6%	42 3.3%	3,828	2.9%
2021年2月7日					26 2.3%	46 3.6%	4,559	3.5%
2021年2月8日					38 3.4%	26 2.1%	4,772	3.6%
2021年2月9日					29 2.6%	25 2.0%	3,847	2.9%
2021年2月10日					31 2.8%	35 2.8%	3,382	2.6%

深度学习的过程-构建模型

1. 假设带有未知参数(parameters)的函数



日期	观看人数	观看人数	观看人数	观看人数	观看人数	观看人数	观看人数
2019-02-25	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-02-26	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-02-27	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-02-28	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-01	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-02	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-03	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-04	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-05	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-06	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-07	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-08	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-09	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-10	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-11	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-12	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-13	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-14	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-15	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-16	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-17	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-18	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-19	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-20	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-21	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-22	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-23	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-24	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-25	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-26	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-27	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-28	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-29	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-30	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48
2019-03-31	19	17,222	26,813	27,602,732	2,656,634	266,778.8	7.48

$$y = f($$



模型 $y = b + wx_1$ 根据先验知识

y : 2月26号的观看人数, x_1 : 2月25号的观看人数 **特征**

w 和 b 是**未知的参数**(从数据中学习得到)

权重 (weight)

偏置 (bias)

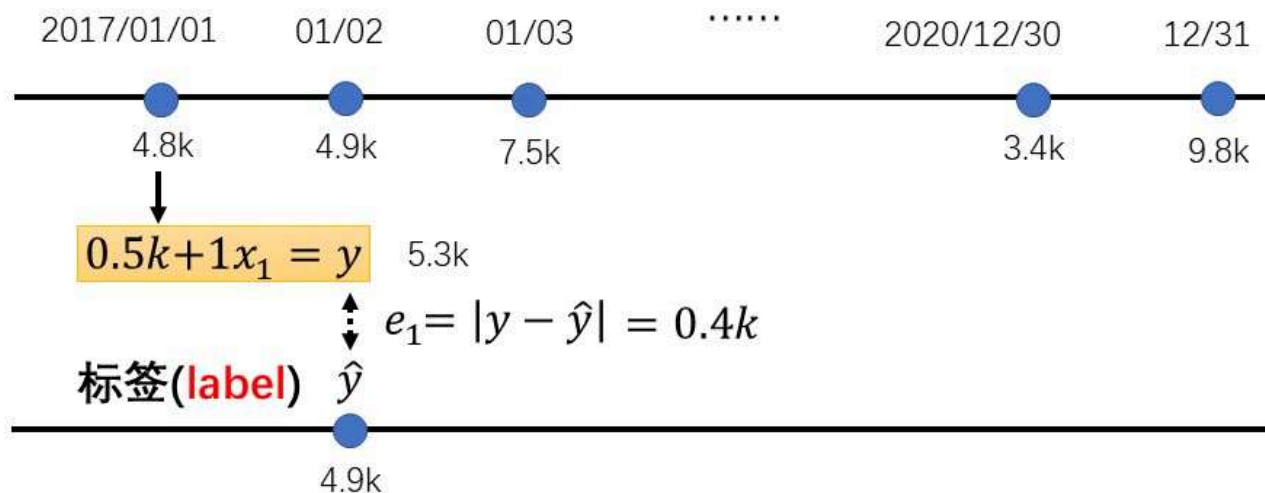
深度学习的过程-定义损失

2. 从训练数据定义损失(Loss)

- Loss 是参数的函数 $L(b, w)$
- Loss: 一组参数有多好的度量标准.

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

从2017/01/01 – 2020/12/31的数据

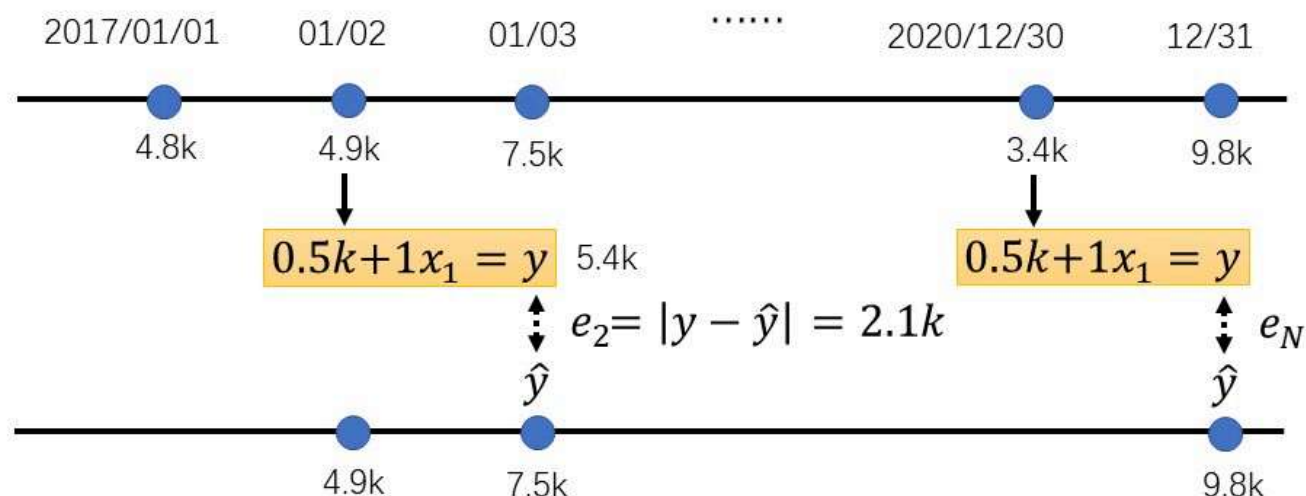


深度学习的过程-定义损失

2. 从训练数据定义损失
- Loss 是参数的函数 $L(b, w)$
 - Loss: 一组参数有多好的度量标准.

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

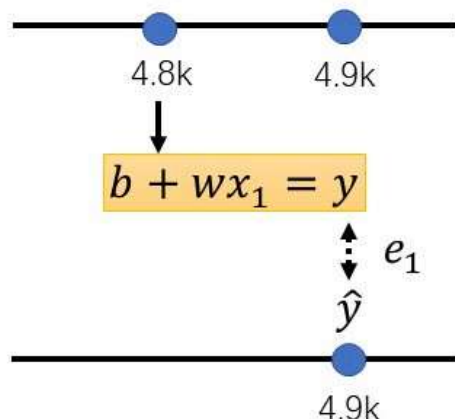
从2017/01/01 – 2020/12/31的数据



深度学习的过程 - 定义损失

2. 从训练数据定义损失(Loss)

- Loss 是参数的函数 $L(b, w)$
- Loss: 一组参数有多好的度量标准.



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

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

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

如果 y 和 \hat{y} 都是几率分布 \longrightarrow 交叉熵 (Cross-entropy)

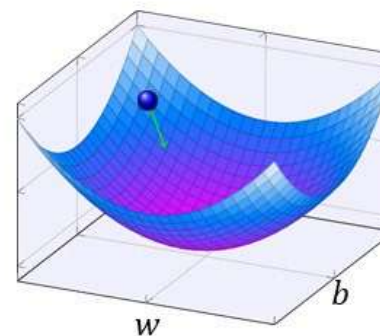
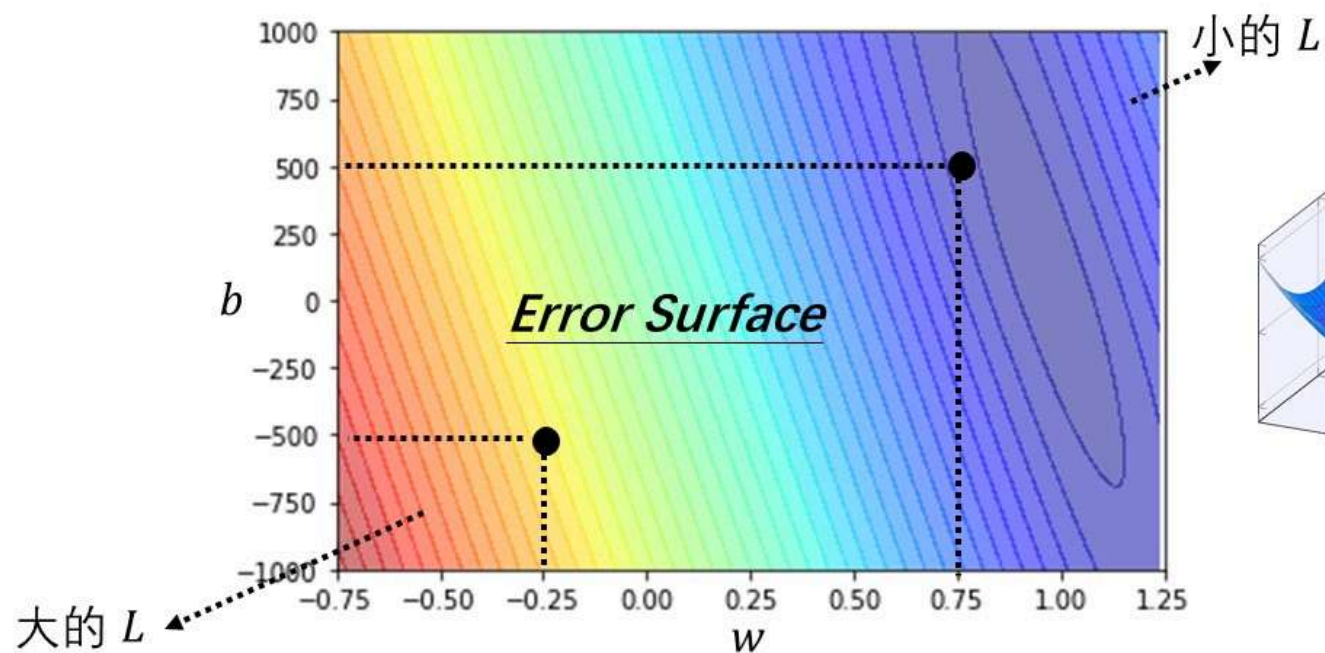


深度学习的过程-定义损失

2. 从训练数据定义损失(Loss)

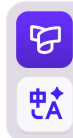
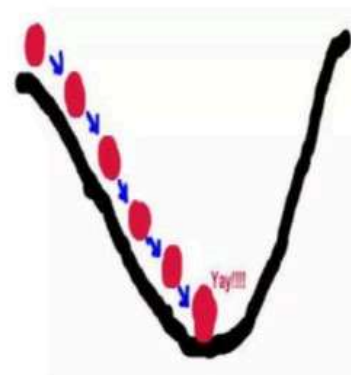
- Loss 是参数的函数 $L(b, w)$
- Loss: 一组参数有多好的度量标准.

模型 $y = b + wx_1$



梯度下降法 (Gradient Descent)

- 假设一个人需要从山的某处开始下山，尽快到达山底。在下山之前他需要确认两件事：下山的方向和下山的距离。因为下山的路有很多，他必须利用一些信息，找到从该处开始最陡峭的方向下山，这样可以保证他尽快到达山底。此外，这座山最陡峭的方向并不是一成不变的，每当走过一段规定的距离，他必须停下来，重新利用现有信息找到新的最陡峭的方向。通过反复进行该过程，最终抵达山底。
- 梯度下降法用于求解无约束最优化问题：山代表了需要优化的函数表达式；山的最低点就是该函数的最优值；每次下山的距离代表后面要解释的学习率；寻找方向利用的信息即为样本数据；最陡峭的下山方向则与函数表达式梯度的方向有关，之所以要寻找最陡峭的方向，是为了满足最快到达山底的限制条件；某处——代表了我们对优化函数设置的初始值，算法后面正是利用这个初始值进行不断的迭代求出最优解。



梯度下降法 (Gradient Descent)

损失函数

$$L(w, b) = \frac{1}{2n} \sum_{i=1}^n [y(x_i) - \hat{y}]^2$$

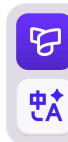
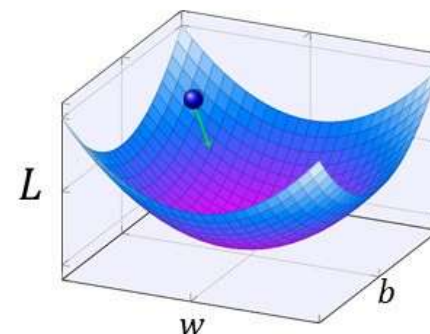
假设要最小化某些函数 $L(\theta)$,它可以是任意的多元实值函数, $\theta = \theta_1; \theta_2; \dots$ 。 θ 代替了 w 和 b 以强调它可能是任意的函数。

当我们在 θ_1 和 θ_2 方向分别将球体移动一个很小的量, 即 $\Delta\theta_1$ 和 $\Delta\theta_2$ 时, 球体将会发生什么情况。微积分告诉我们 L 将会有如下变化:

$$\Delta L = \frac{\partial L}{\partial \theta_1} \Delta\theta_1 + \frac{\partial L}{\partial \theta_2} \Delta\theta_2 \quad \Delta L \text{ 符号?}$$

要寻找一种选择 $\Delta\theta_1$ 和 $\Delta\theta_2$ 的方法使得 ΔL 为负; 选择负是为了让球体滚落。

$$\Delta\theta = (\Delta\theta_1, \Delta\theta_2)^T, \quad \nabla L = \left(\frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2} \right) \quad \longrightarrow \quad \Delta L = \nabla L \cdot \Delta\theta$$



梯度下降法 (Gradient Descent)

$$\Delta L = \nabla L \cdot \Delta \theta$$

这个方程真正让我们兴奋的是它让我们看到了如何选取 $\Delta \theta$ 才能让 ΔL 为负数。
假设我们选取：

$$\Delta \theta = -\eta \nabla L$$

η 称为学习速率。

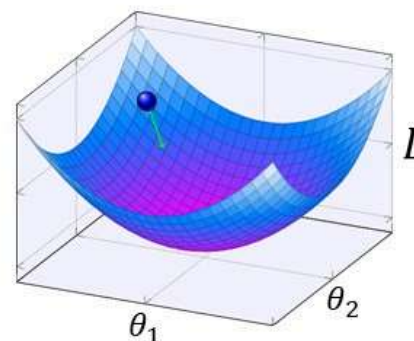
$$\Delta L = \nabla L \cdot \Delta \theta = -\eta \nabla L \cdot \nabla L = -\eta \|\nabla L\|^2$$

$$\theta \rightarrow \theta' = \theta - \eta \nabla L$$

然后我们用它再次更新规则来计算下一次移动。如果我们反复持续这样做，我们将持续减小 L 直到——正如我们希望的——获得一个全局的最小值。

总结一下，梯度下降算法工作的方式就是重复计算梯度 ∇L ，然后沿着相反的方向移动，沿着山谷“滚落”。我们可以想象它像这样：

$$\nabla L = \left(\frac{\partial L}{\partial \theta_1}, \dots, \frac{\partial L}{\partial \theta_m} \right) \quad \Delta \theta = -\eta \nabla L$$

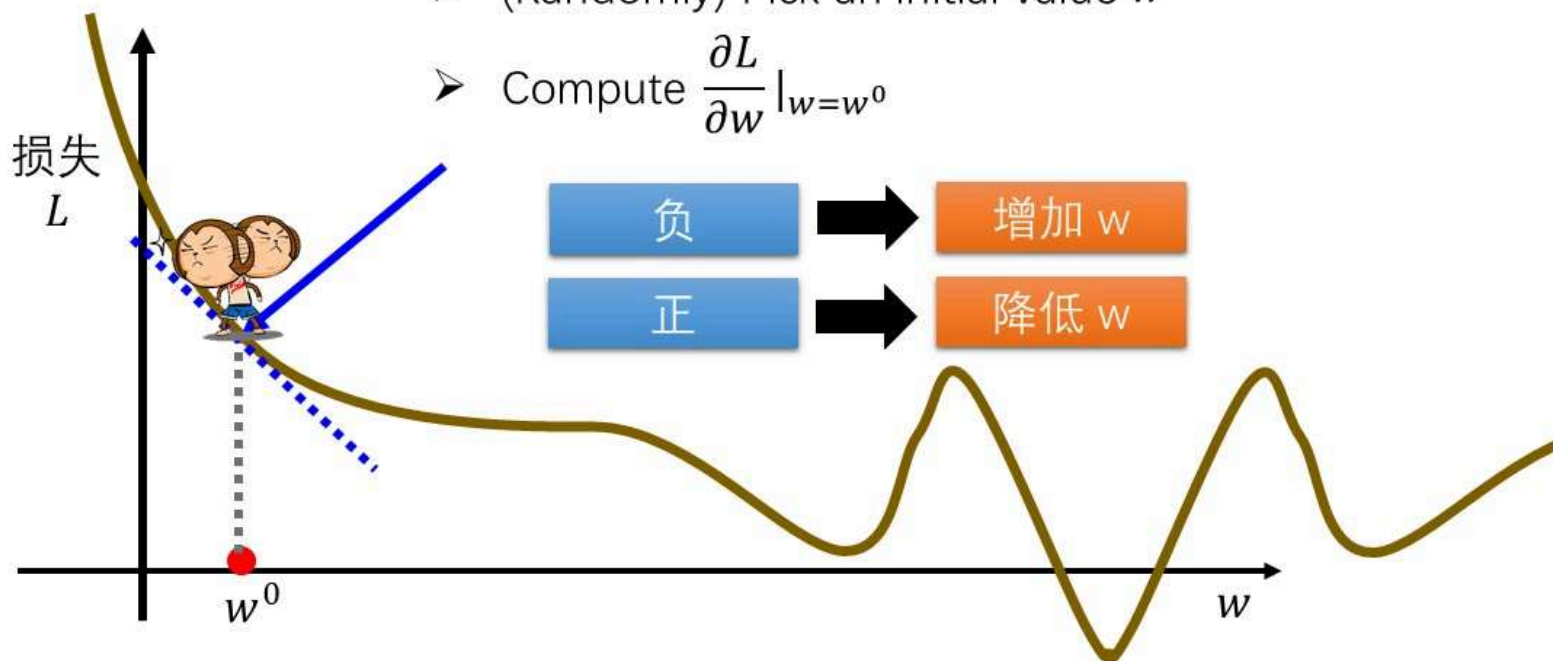


3. 优化(Optimization)

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

梯度下降 (Gradient Descent)

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



3. 优化(Optimization)

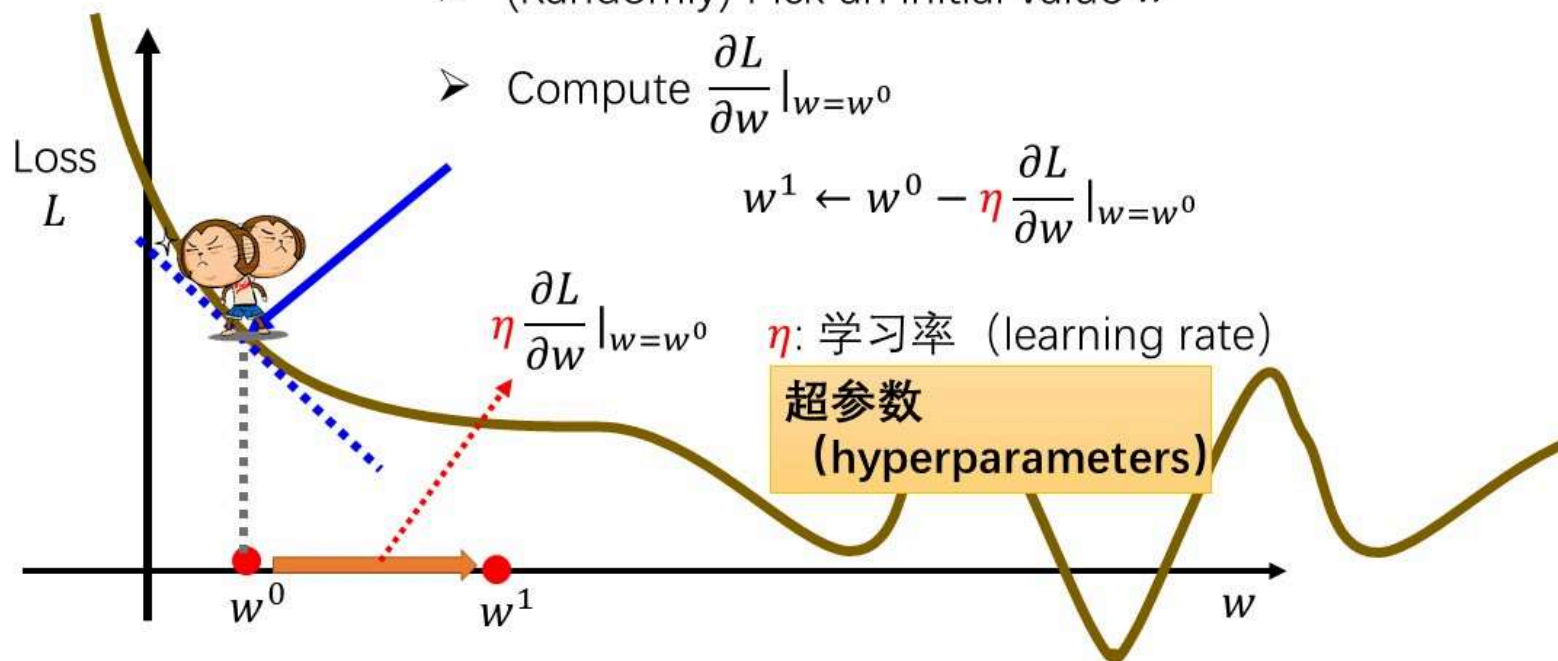
$$w^* = \arg \min_w L$$

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. 优化Optimization

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

Gradient Descent

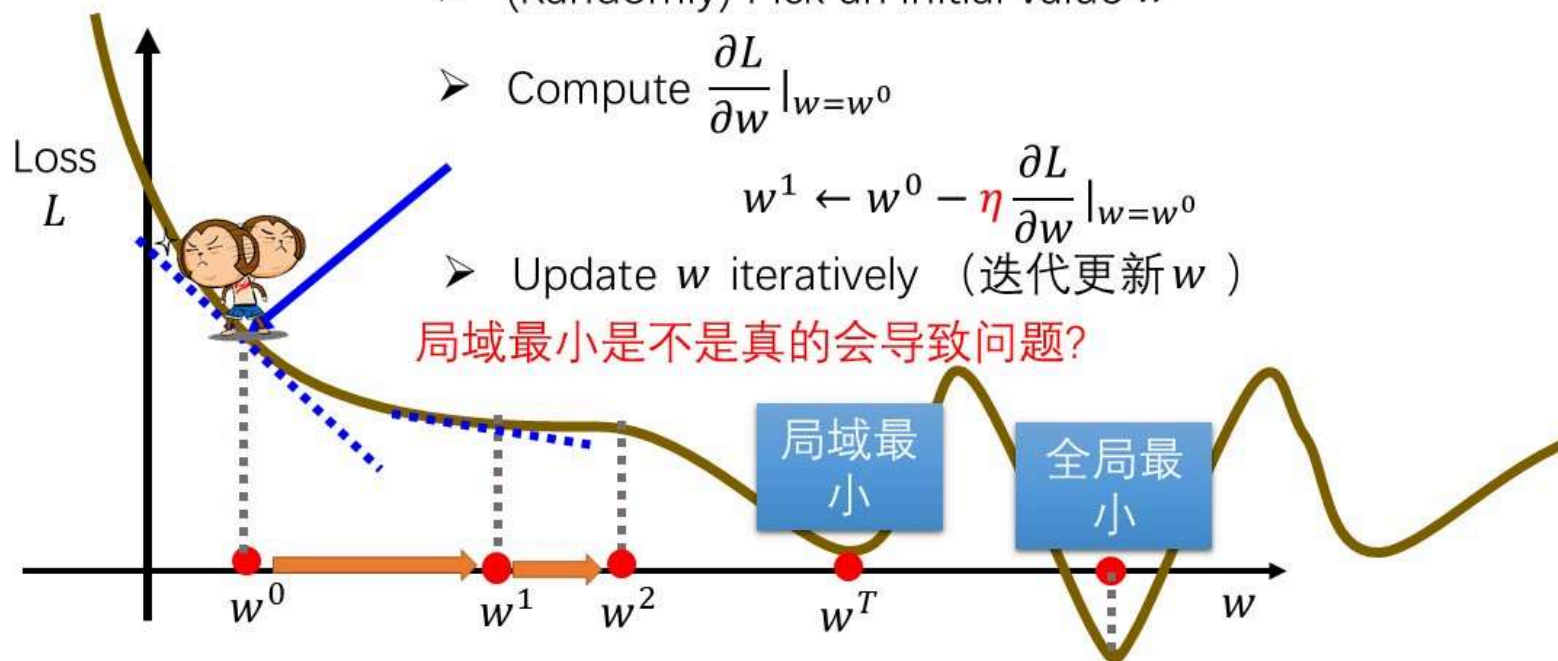
➤ (Randomly) Pick an initial value w^0

➤ 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 (迭代更新 w)

局域最小是不是真的会导致问题?



深度学习的过程-优化

3. Optimization

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

➤ 随机初始化 w^0, b^0

➤ 计算

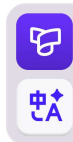
$$\begin{aligned} \frac{\partial L}{\partial w} \big|_{w=w^0, b=b^0} \\ \frac{\partial L}{\partial b} \big|_{w=w^0, b=b^0} \end{aligned}$$

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

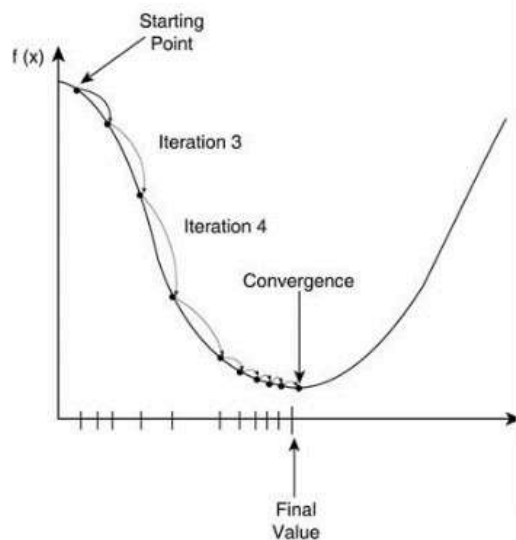
$$b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b} \big|_{w=w^0, b=b^0}$$

Can be done in one line in most deep learning frameworks (一行代码解决)

➤ Update (更新) w and b iteratively (迭代)



梯度下降法 (Gradient Descent)



搜索中步长 η 也叫作**学习率**
(Learning Rate)

1. 给定待优化连续可微函数 $L(\Theta)$ 、学习率 η 以及一组初始值 $L_0 = (\theta_{01}, \theta_{02}, \dots, \theta_{0l},)$
2. 计算待优化函数梯度: $\nabla L(\Theta_0)$
3. 更新迭代公式: $\Theta^{0+1} = \Theta_0 - \eta \nabla L(\Theta_0)$
4. 计算 Θ^{0+1} 处函数梯度 $\nabla L(\Theta_{0+1})$
5. 计算梯度向量的模来判断算法是否收敛: $\|\nabla L(\Theta)\| \leq \varepsilon$
6. 若收敛, 算法停止, 否则根据迭代公式继续迭代

➤ 经过迭代计算风险函数的最小值

$$\theta_{t+1} = \theta_t - \eta \frac{\partial L}{\partial \theta}$$

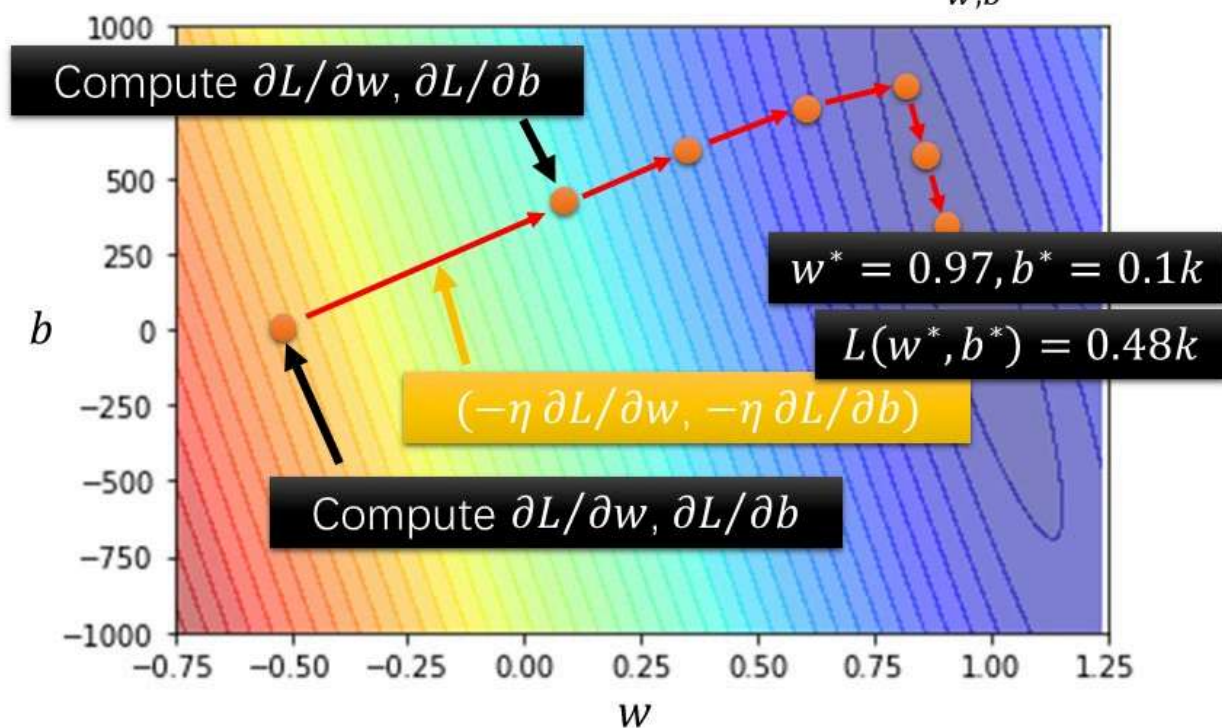


深度学习的过程-优化

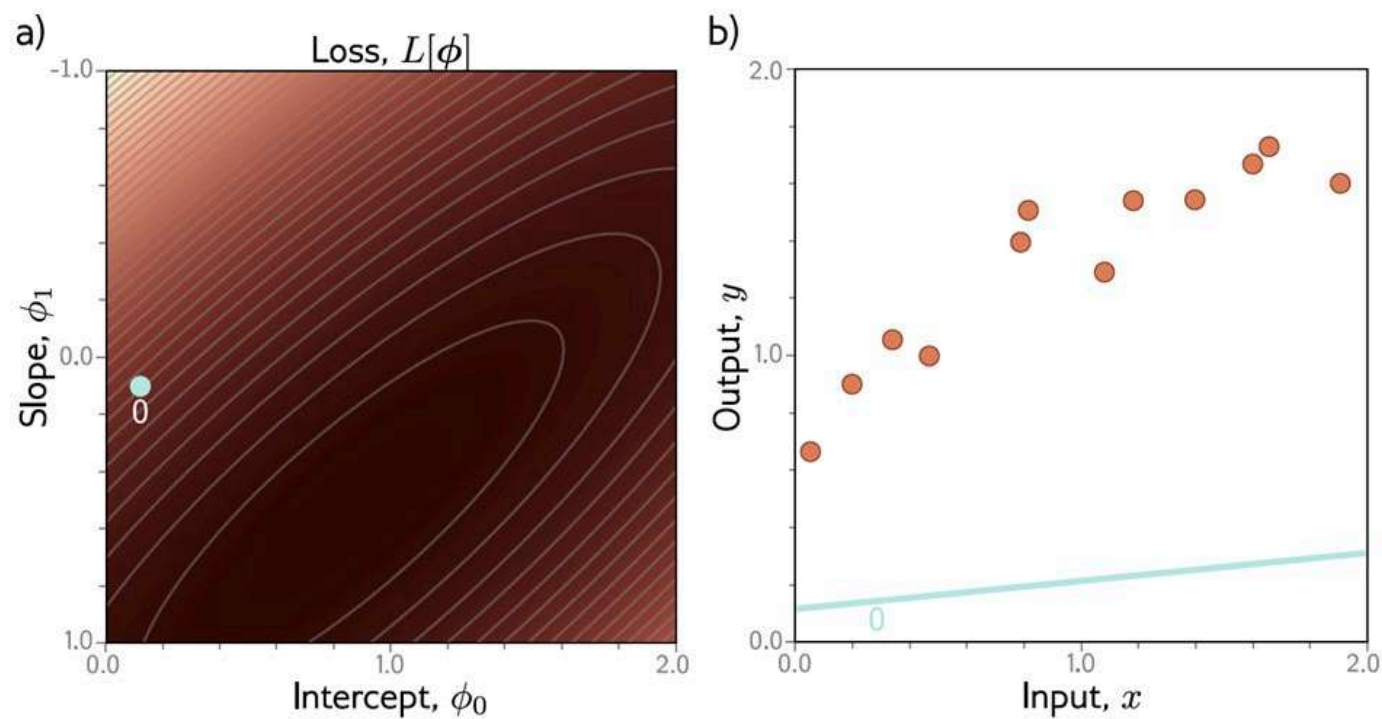
3. Optimization

模型 $y = b + wx_1$

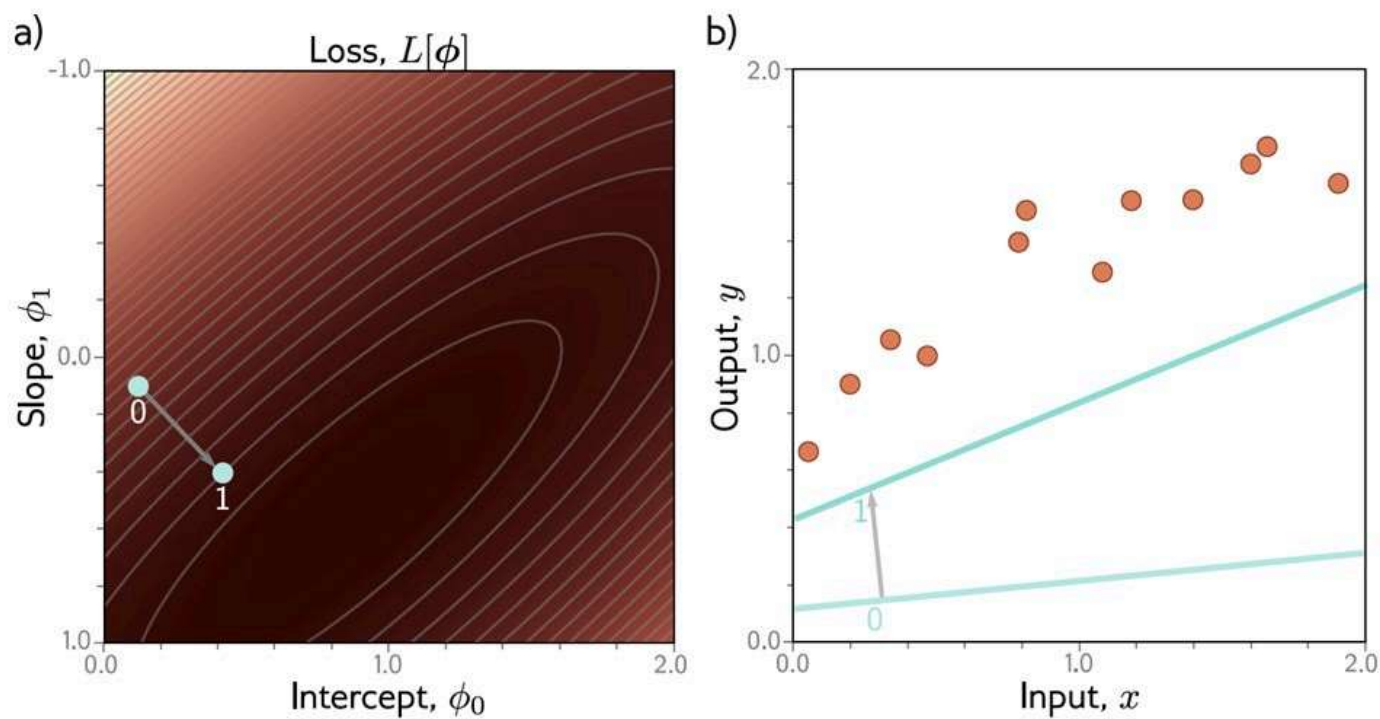
$$w^*, b^* = \arg \min_{w, b} L$$



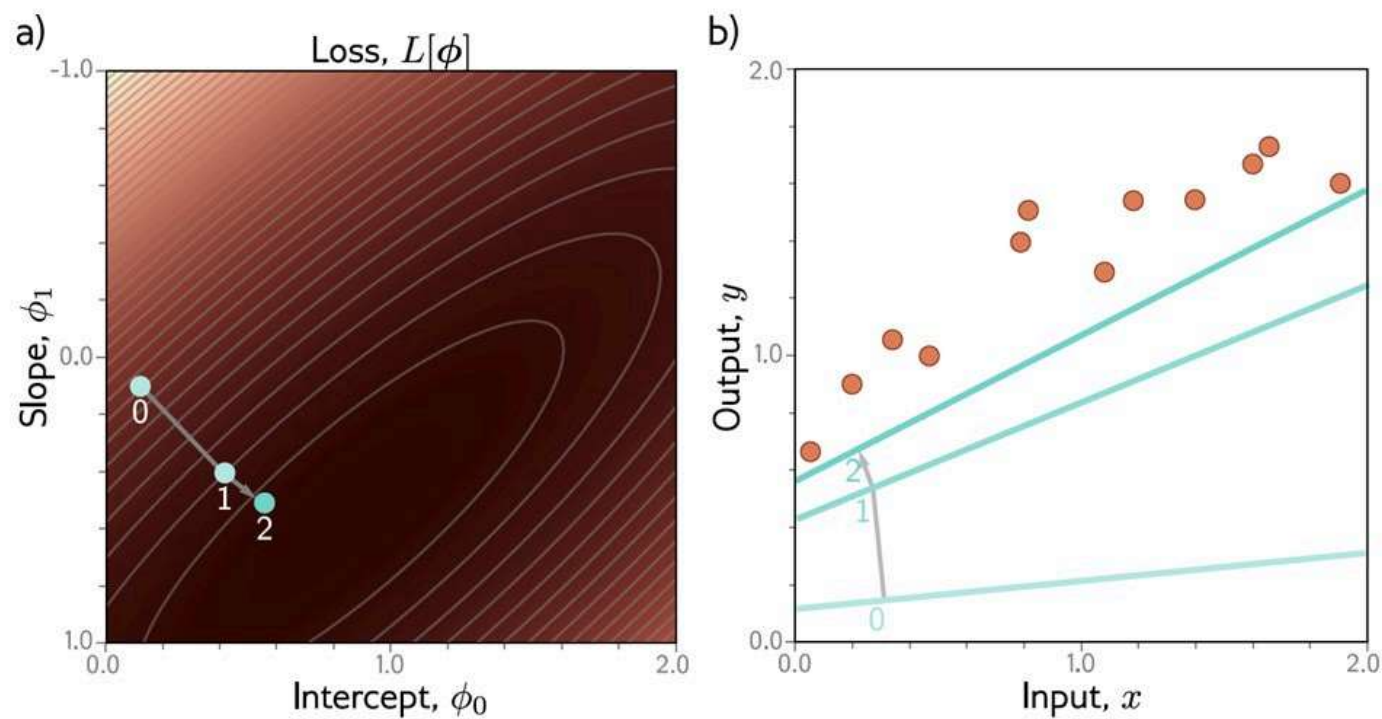
Example: 1D Linear regression training



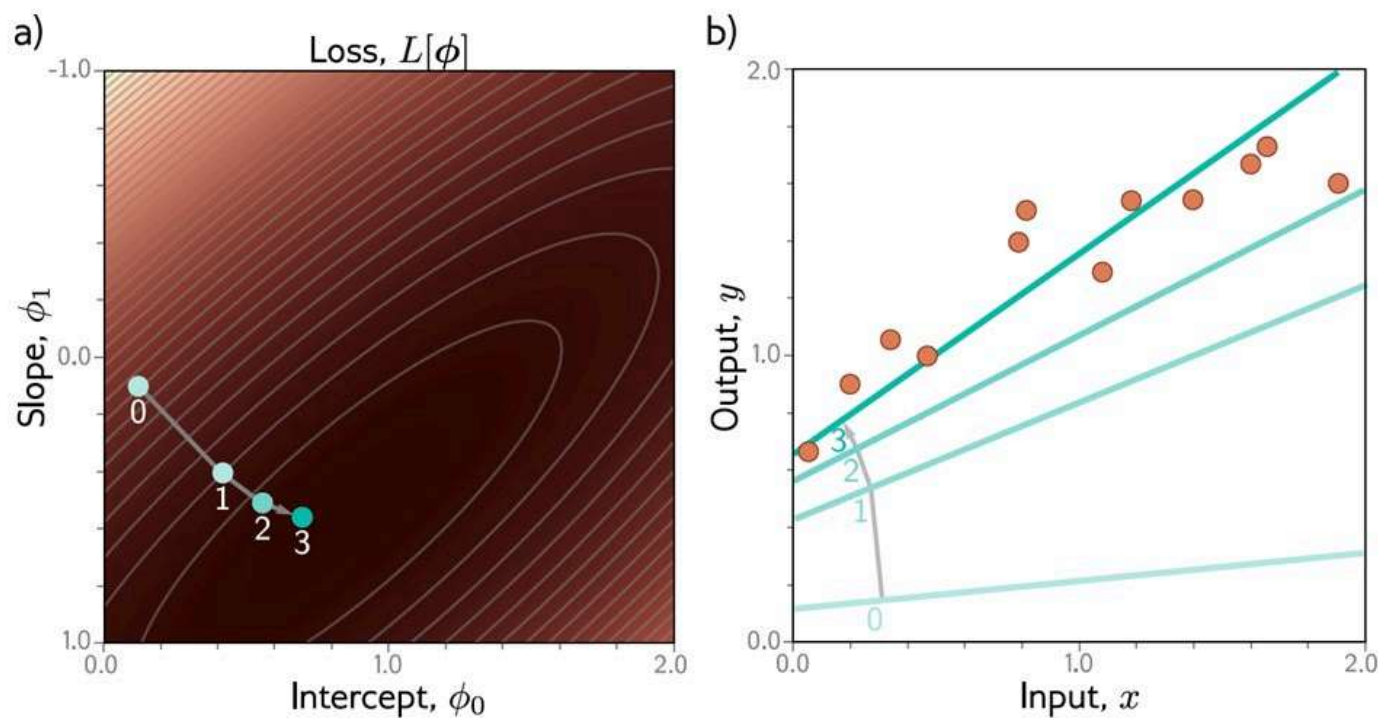
Example: 1D Linear regression training



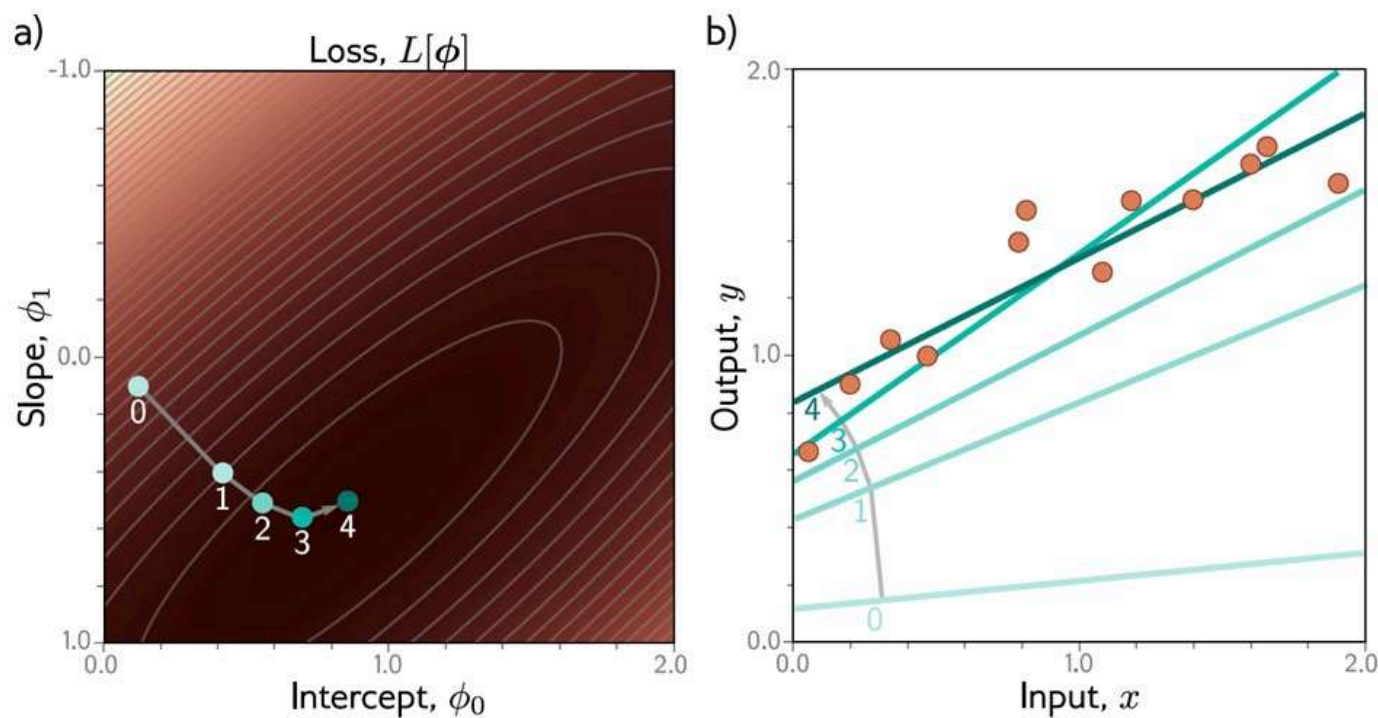
Example: 1D Linear regression training



Example: 1D Linear regression training



Example: 1D Linear regression training



This technique is known as **gradient descent**

深度学习的过程-构建模型-定义损失-优化

deep Learning is so simple



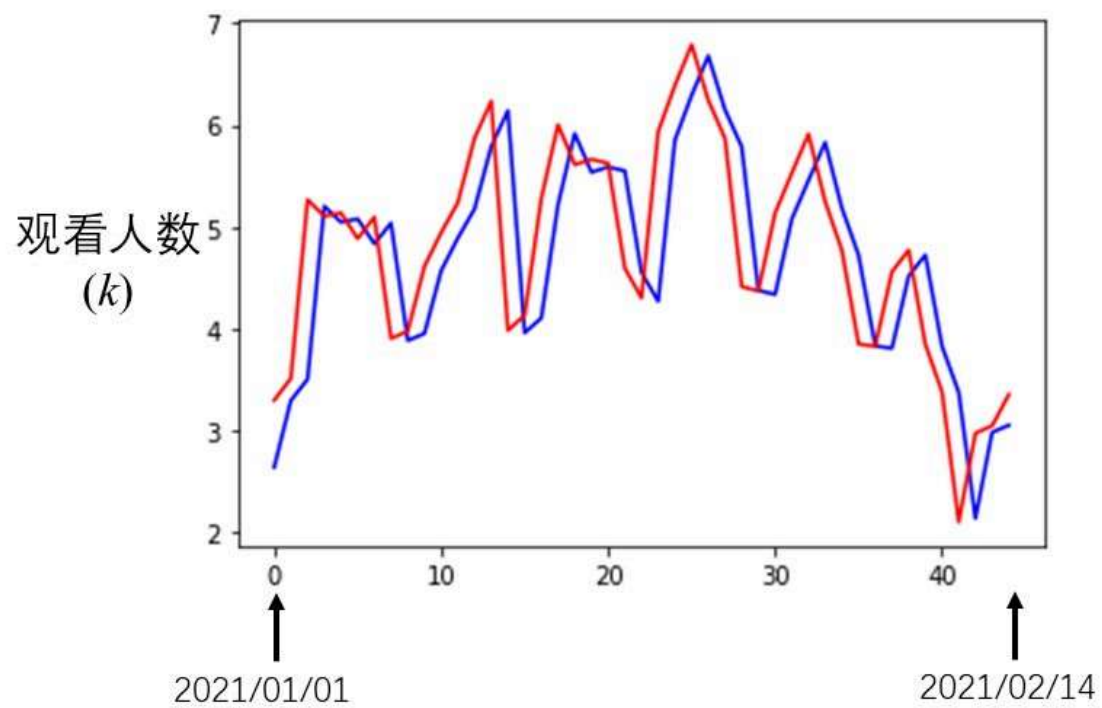
$y = 0.1k + 0.97x_1$ 得到最小的损失 $L = 0.48k$
从数据2017 – 2020 (训练数据 **training data**)

How about data of 2021 (测试数据 **unseen during training**)?

$$L' = 0.58k$$

$y = 0.1k + 0.97x_1$

Red: 后台统计数据
blue: 模型预测数据



$$y = b + wx_1 \quad \begin{array}{cc} 2017 - 2020 & 2021 \\ L = 0.48k & L' = 0.58k \end{array}$$

$$y = b + \sum_{j=1}^7 w_j x_j \quad \begin{array}{cc} 2017 - 2020 & 2021 \\ L = 0.38k & L' = 0.49k \end{array}$$

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

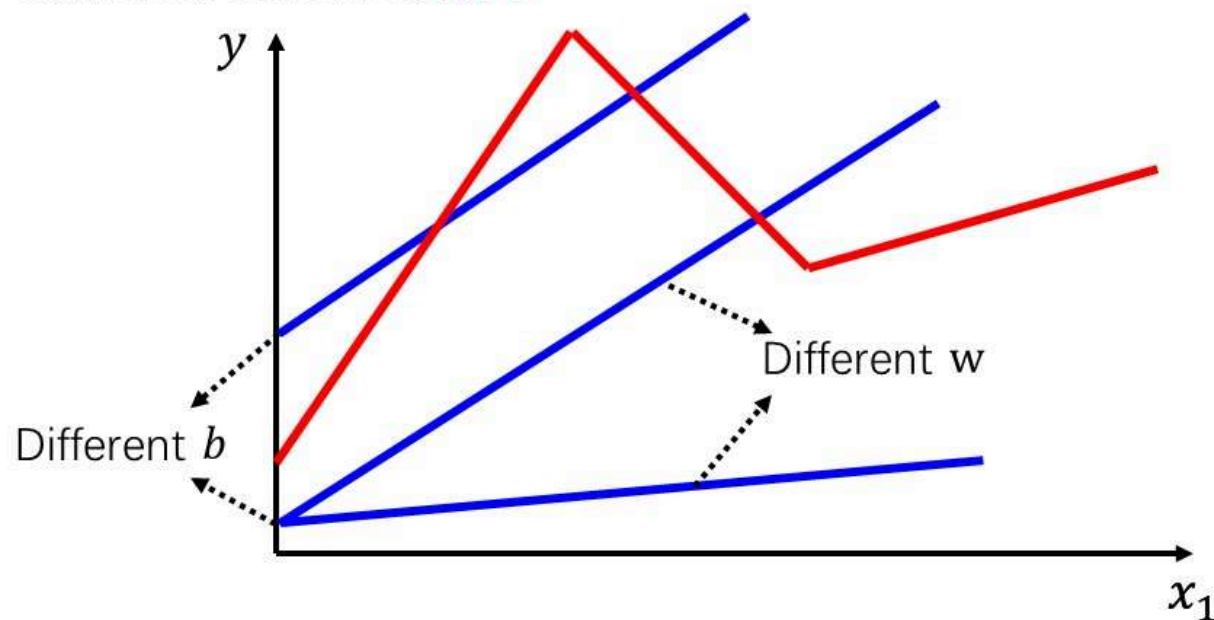
$$y = b + \sum_{j=1}^{28} w_j x_j \quad \begin{array}{cc} 2017 - 2020 & 2021 \\ L = 0.33k & L' = 0.46k \end{array}$$

$$y = b + \sum_{j=1}^{56} w_j x_j \quad \begin{array}{cc} 2017 - 2020 & 2021 \\ L = 0.32k & L' = 0.46k \end{array}$$

线性模型(Linear models)

深度学习非线性模型初探

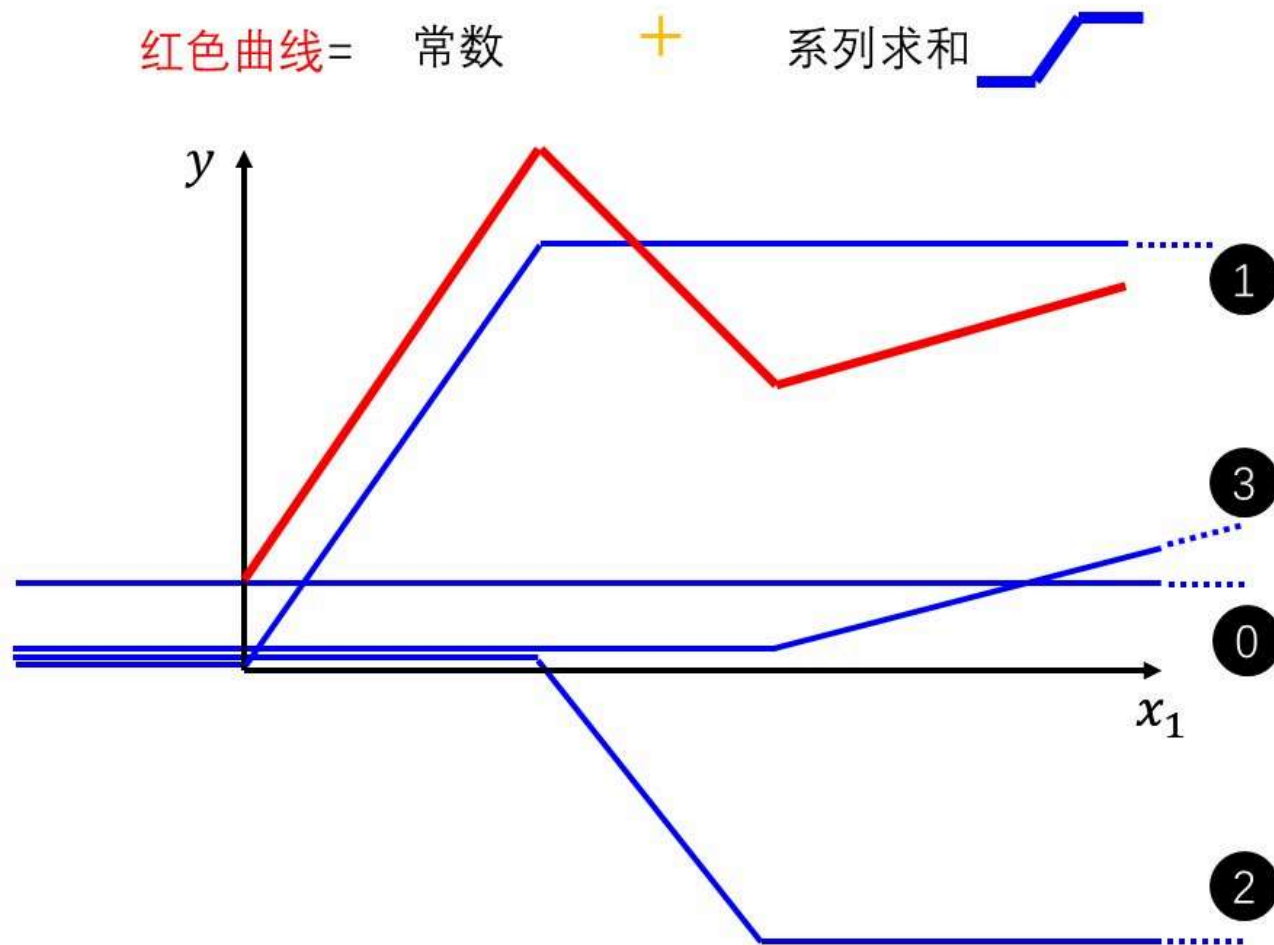
线性模型…我们需要更加复杂的模式.



线性模型具有严重的限制. **模型偏离**

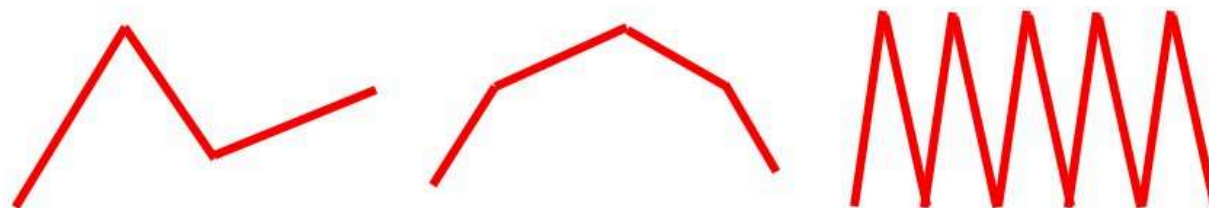
我们需要更加灵活的模型!

红色曲线 = 常数 + 系列求和



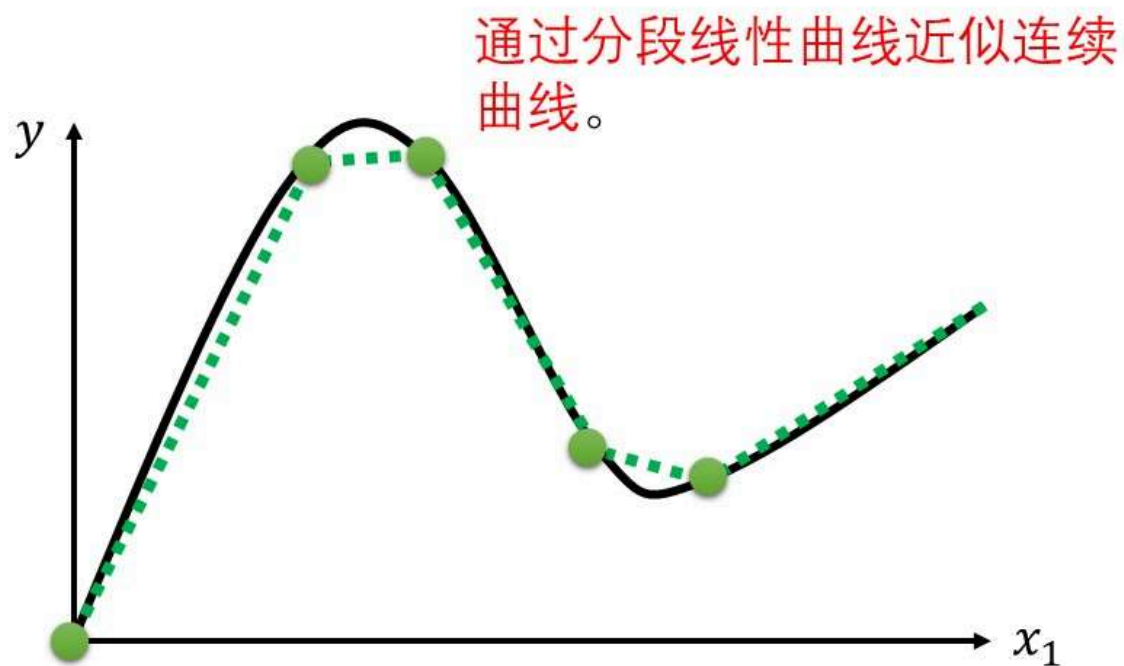
所有分段连续线性曲线

= constant + sum of a set of 



More pieces require more 

超越分段线性连续的的曲线?



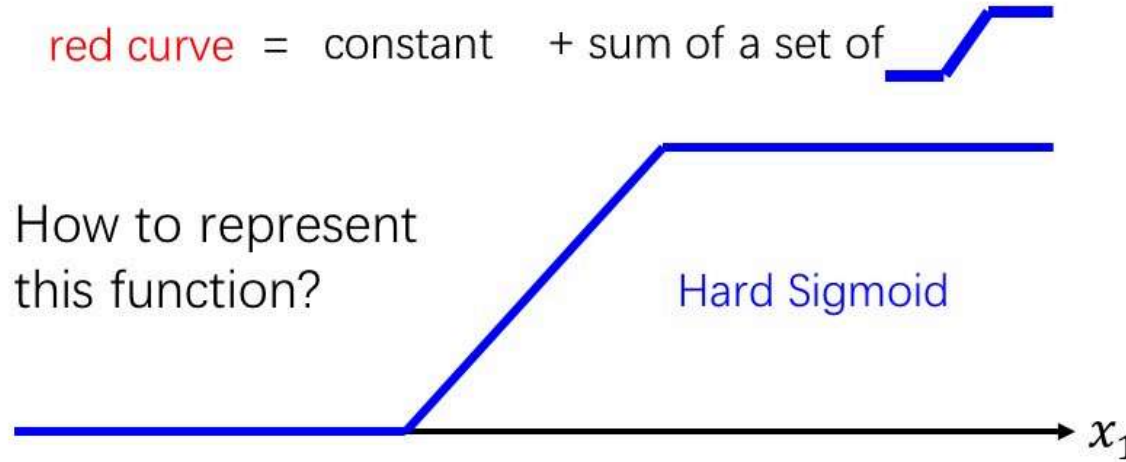
为了获得良好的近似值，我们需要足够的分段。

red curve = constant + sum of a set of



How to represent
this function?

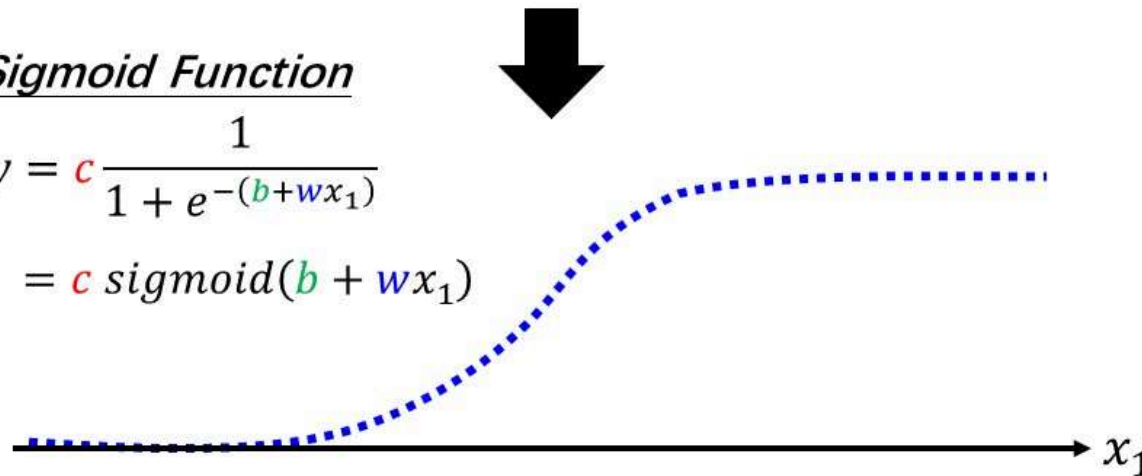
Hard Sigmoid

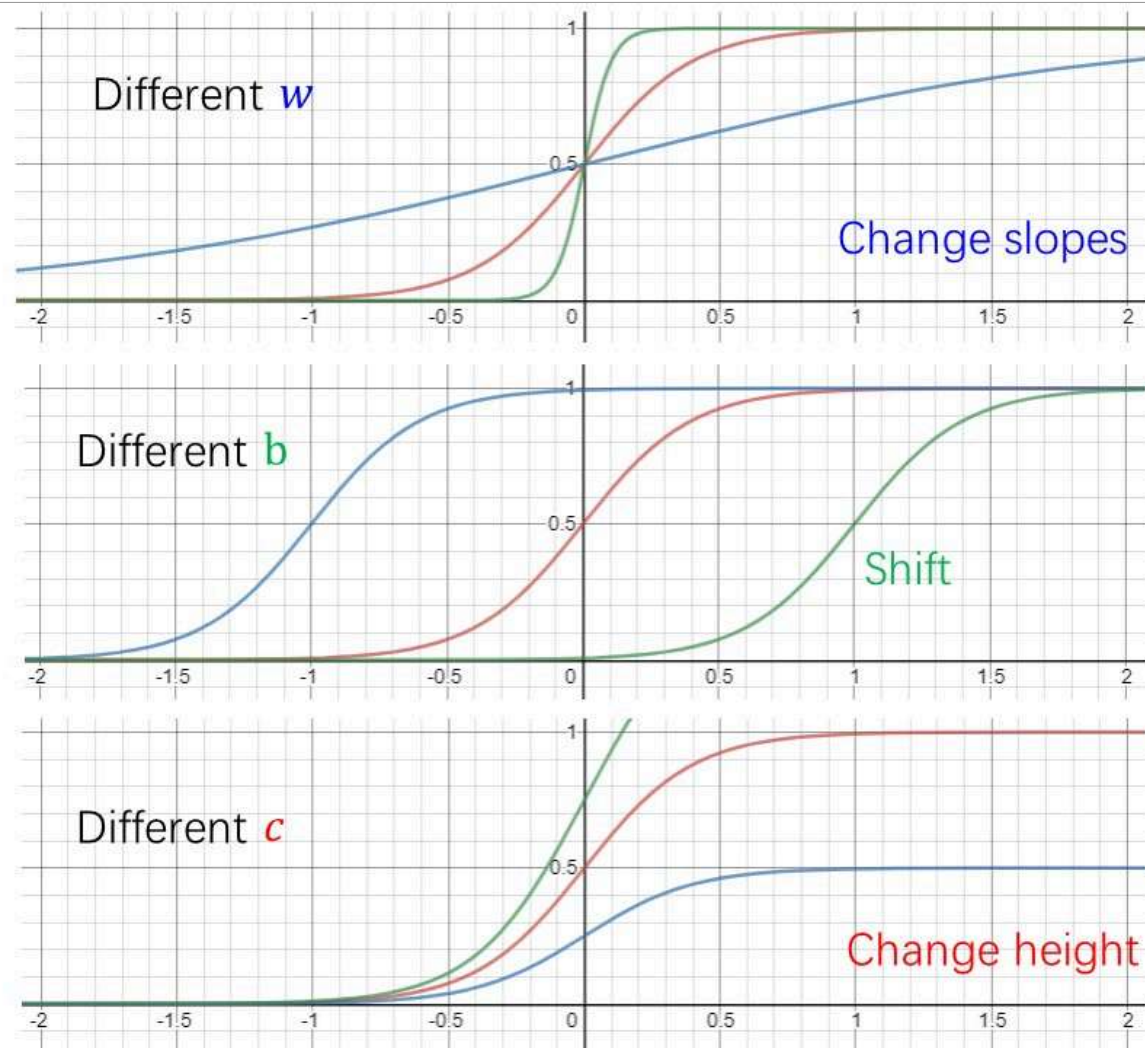


Sigmoid Function

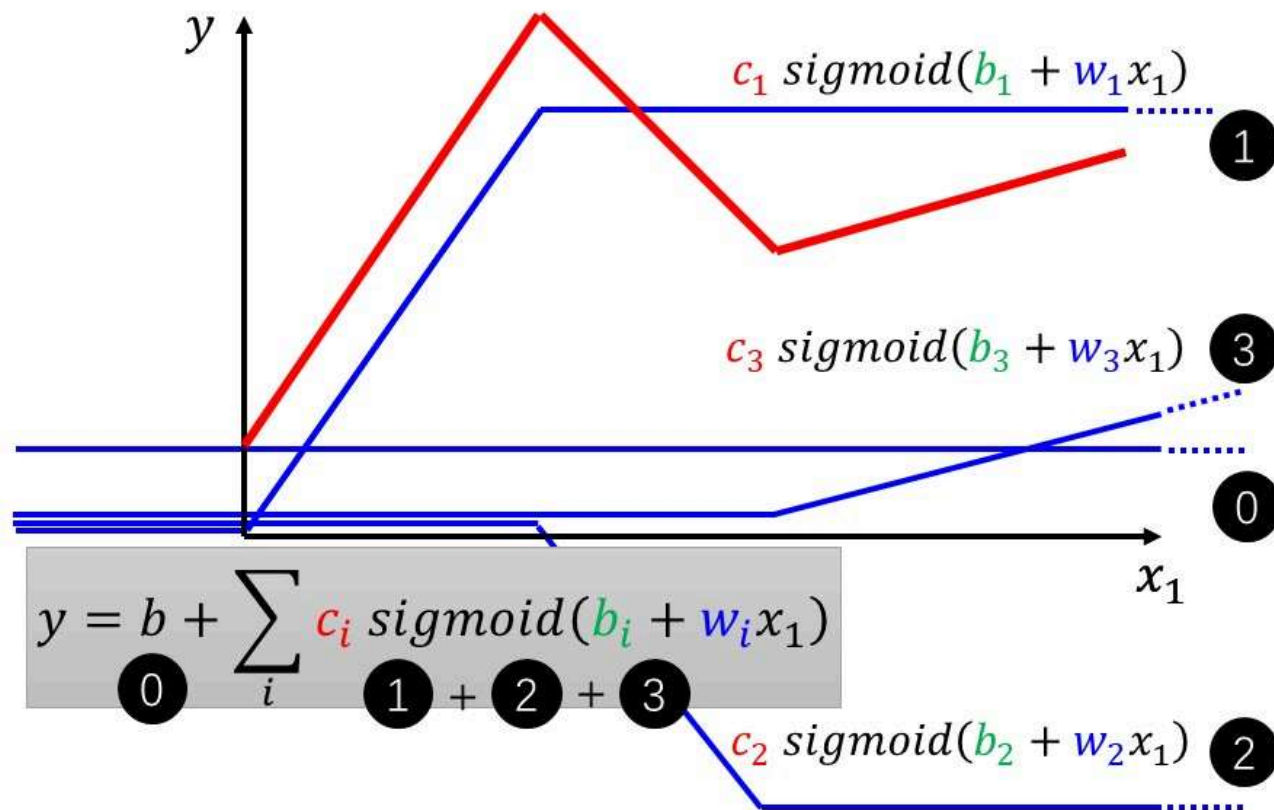
$$y = c \frac{1}{1 + e^{-(b + wx_1)}}$$

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





red curve = sum of a set of  + constant



新模型: 较多特征(More Features)

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

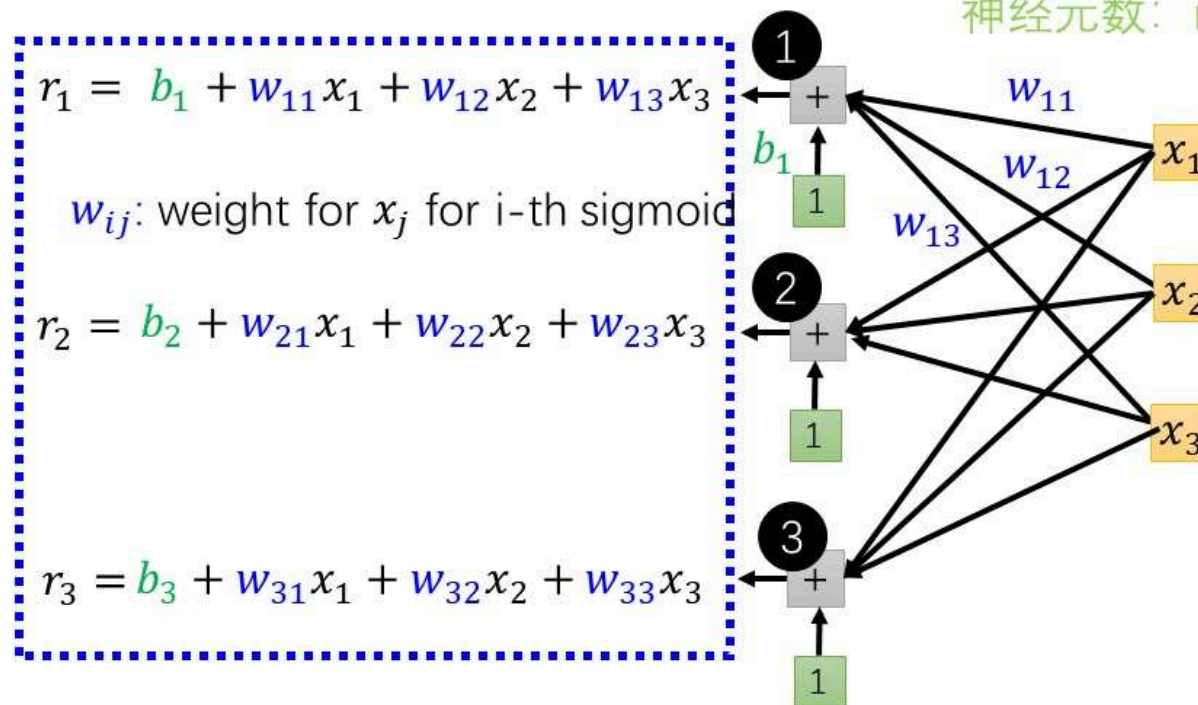
$$y = b + \sum_i c_i \operatorname{sigmoid}(\underline{b_i + w_i x_1})$$

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

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

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

$j: 1, 2, 3$
 特征数: no. of features
 $i: 1, 2, 3$
 神经元数: no. of sigmoid



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

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

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

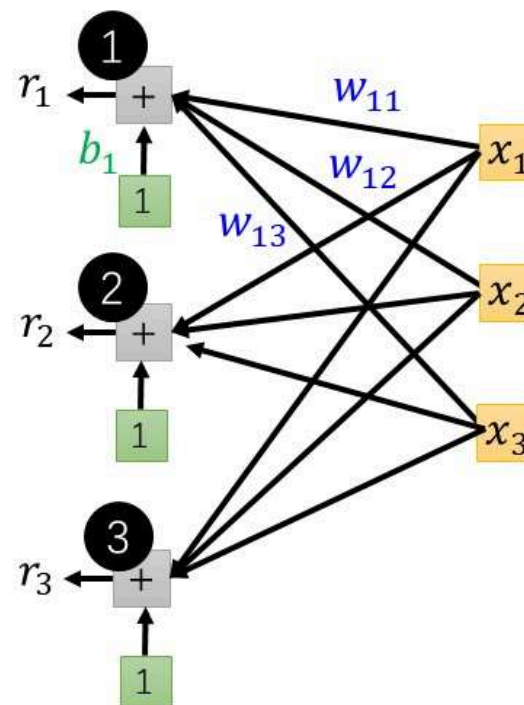
$$r_3 = b_3 + w_{31}x_1 + w_{32}x_2 + w_{33}x_3$$

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

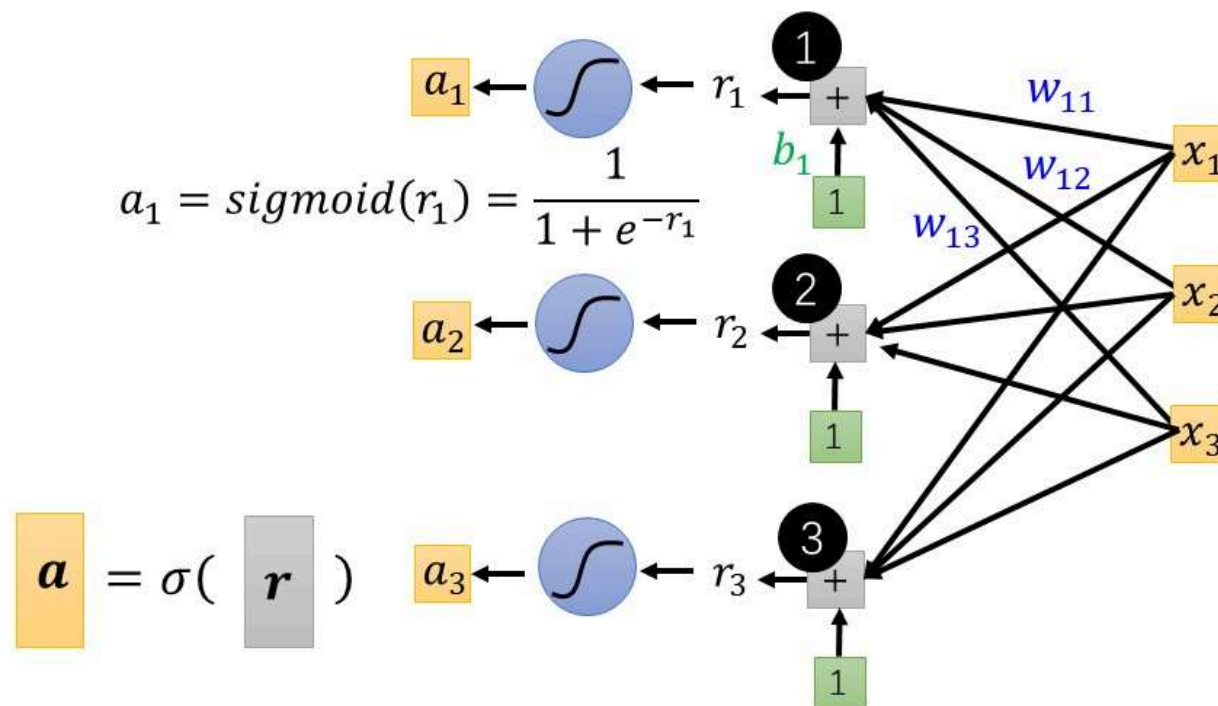
$$\mathbf{r} = \mathbf{b} + \mathbf{W} \mathbf{x}$$

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

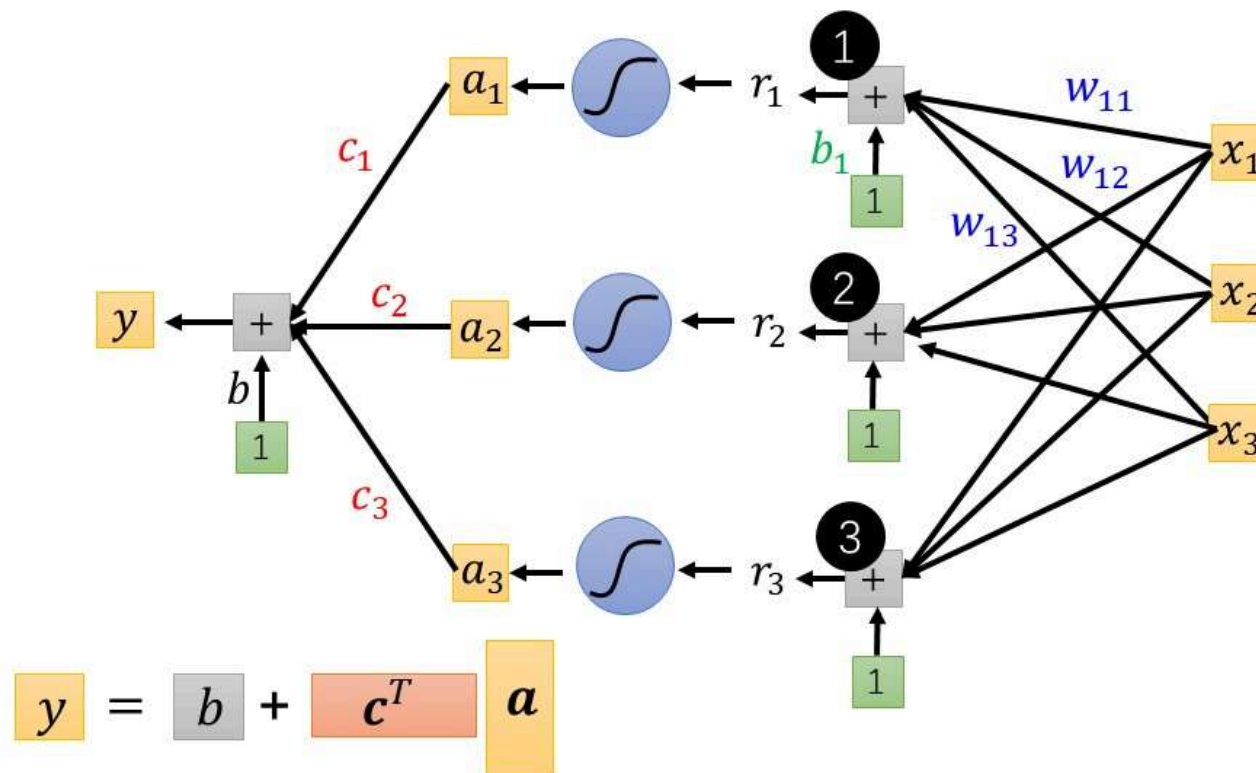
$$\mathbf{r} = \mathbf{b} + \mathbf{W} \mathbf{x}$$

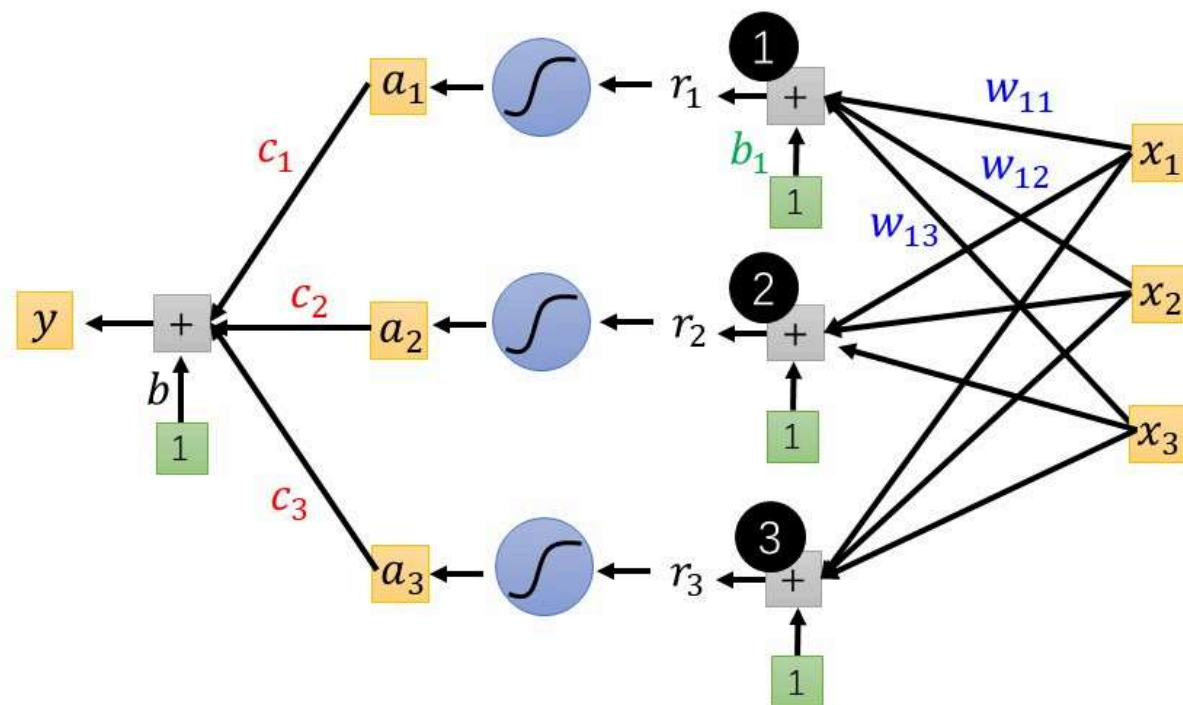


$$y = b + \sum_i c_i \text{sigmoid}\left(b_i + \sum_j w_{ij} x_j\right) \quad \begin{matrix} i: 1,2,3 \\ j: 1,2,3 \end{matrix}$$



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

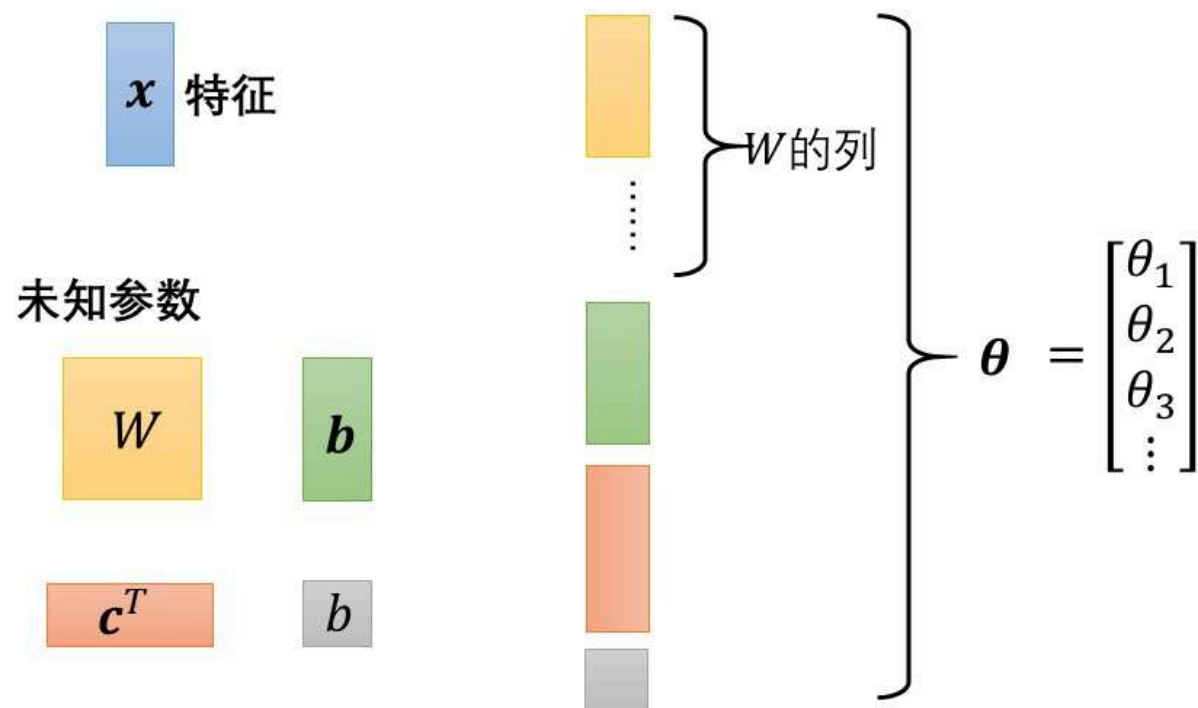




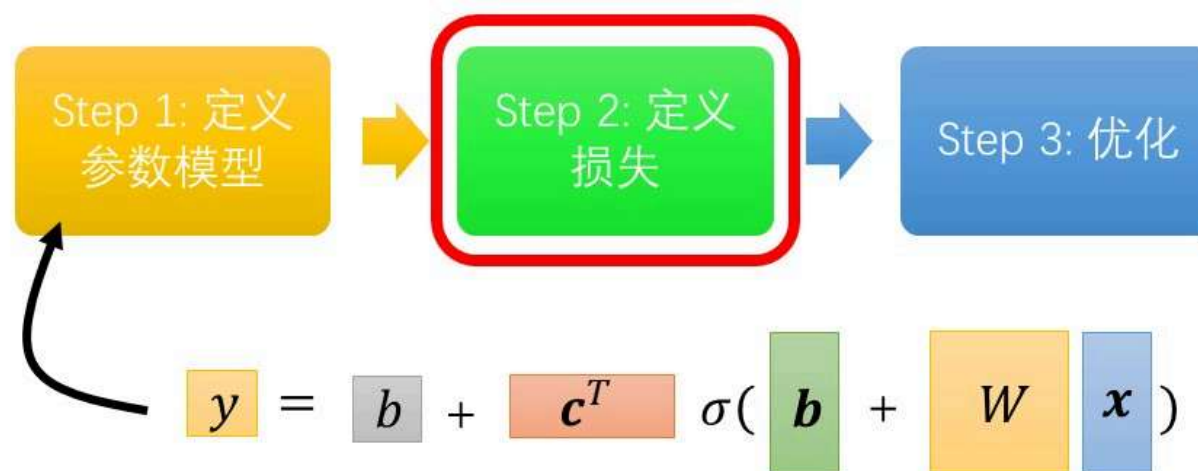
$$y = b + c^T \sigma(b + Wx)$$

带有未知参数的函数

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



回到 DL 框架



损失



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

标签: Label \hat{y}

特征: feature x

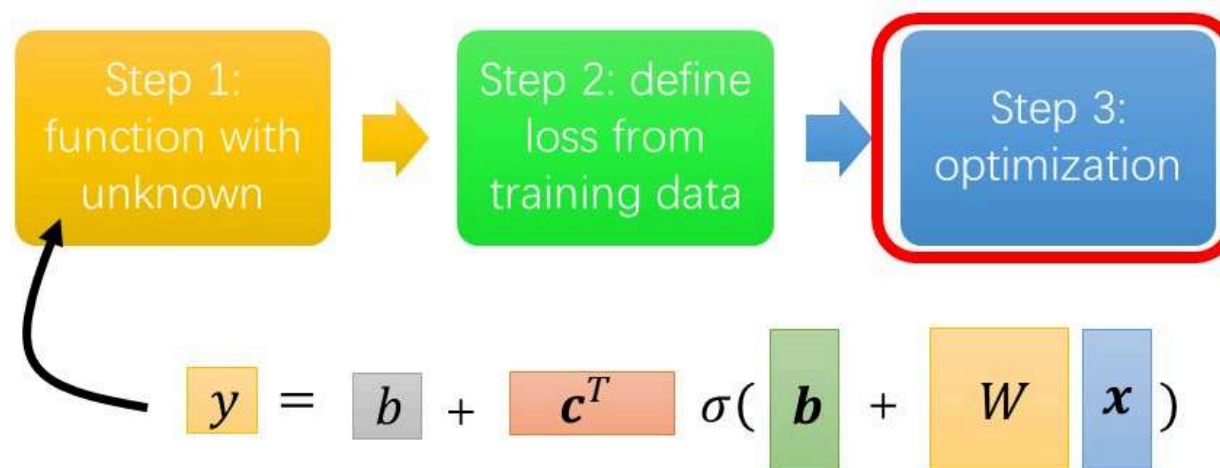
给定一套参数值

The diagram illustrates the model equation $y = b + c^T \sigma(b + Wx)$. It shows the relationship between the target y and the predicted value \hat{y} with an error e . The parameters b , c^T , b , and W are grouped under the label "给定一套参数值" (Given a set of parameter values). The input x is labeled "特征: feature" (Feature: feature).

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

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

Back to DL Framework-优化



新模型的优化: Optimization of New Model

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

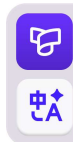
$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \end{bmatrix}$$

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

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

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

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



新模型的优化: Optimization of New Model

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

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

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

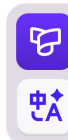
$$\boldsymbol{\theta}^1 \leftarrow \boldsymbol{\theta}^0 - \eta \mathbf{g}$$

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

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

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

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



新模型的优化: Optimization of New Model

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

➤ (Randomly) Pick initial values θ^0

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

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

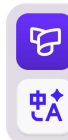
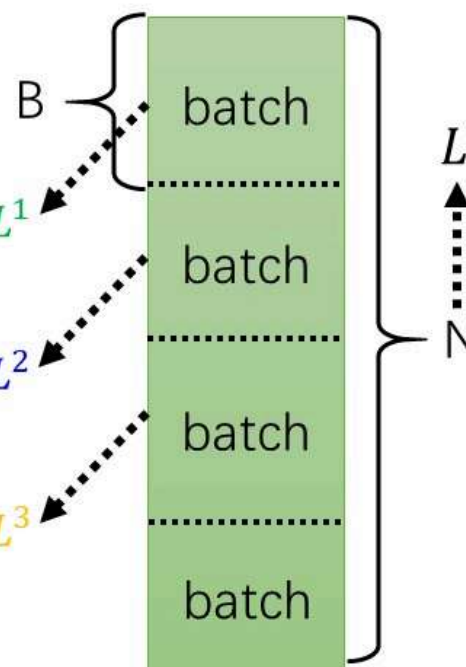
➤ Compute gradient $\mathbf{g} = \nabla L^2(\theta^1)$

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

➤ 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



新模型的优化: Optimization of New Model

例子 1

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

How many update in **1 epoch**?

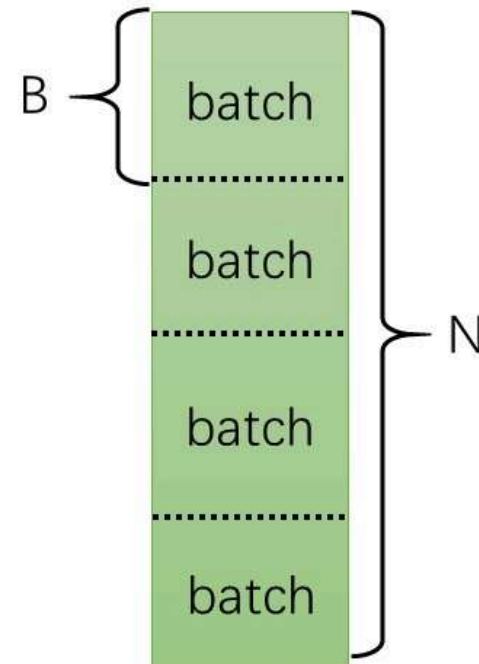
1,000 更新

例子 2

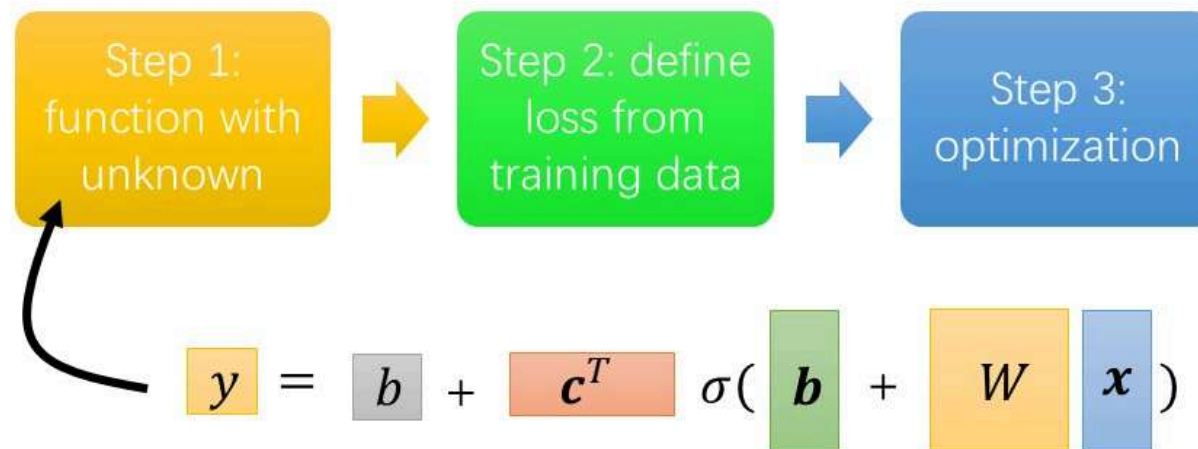
- 1,0000 examples ($N = 1,0000$)
- Batch size is 100 ($B = 100$)

How many update in **1 epoch**?

100updates

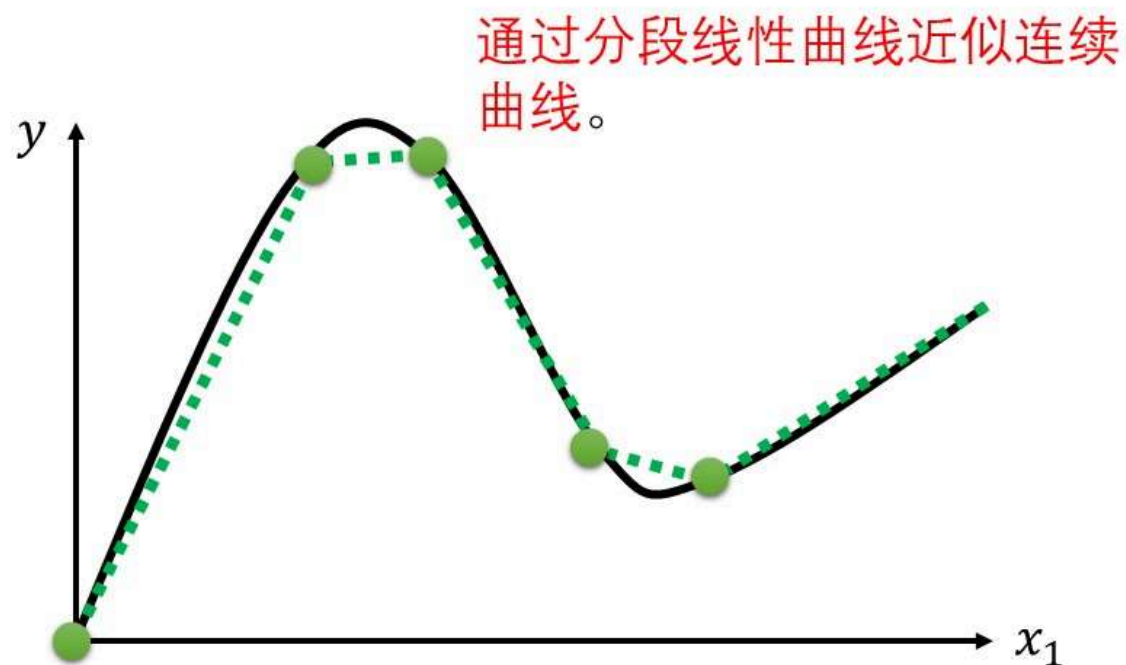


回到深度学习框架



很多的变形模型: More variety of models ...

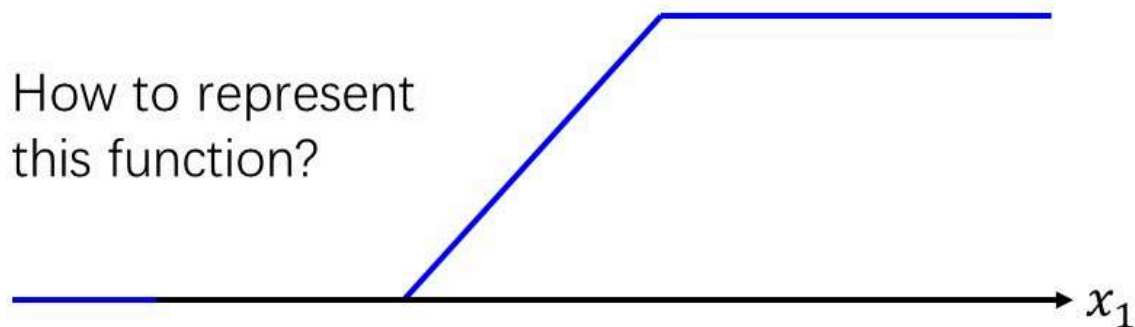
超越分段线性连续的的曲线?



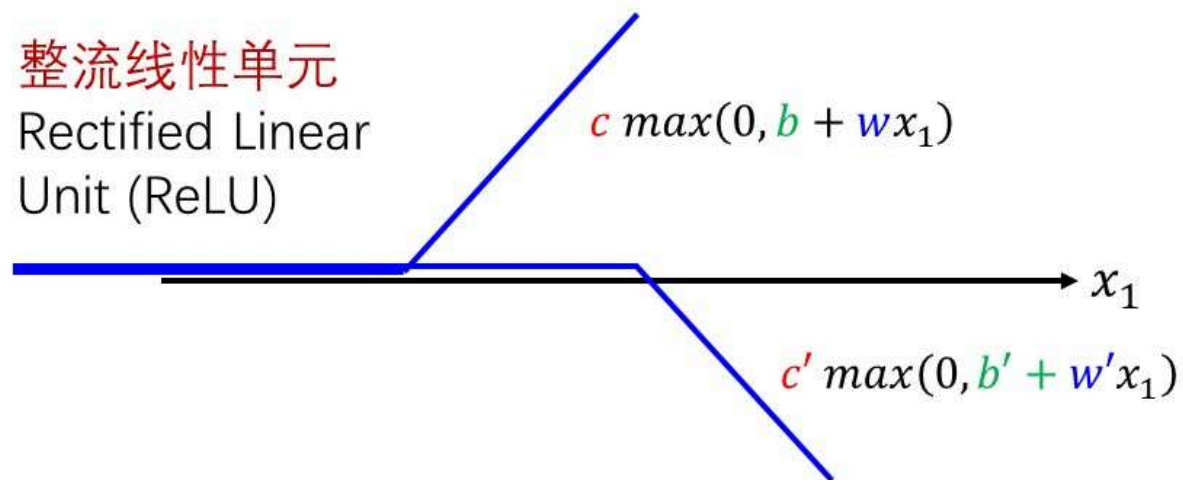
为了获得良好的近似值，我们需要足够的分段。

Sigmoid \rightarrow ReLU

How to represent
this function?



整流线性单元
Rectified Linear
Unit (ReLU)



$$c \max(0, b + wx_1)$$

$$c' \max(0, b' + w'x_1)$$

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

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

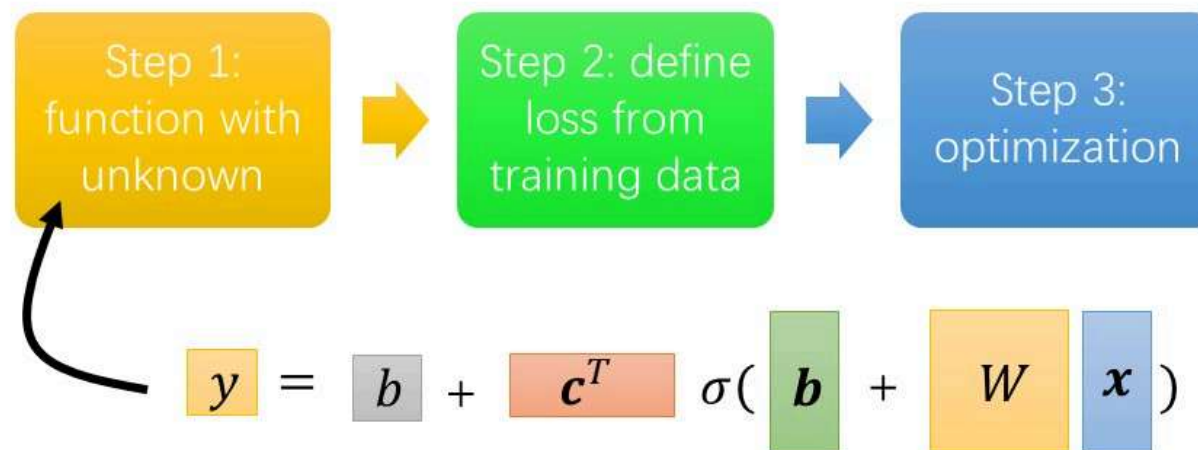
Which one is better?

实验结果:Experimental Results

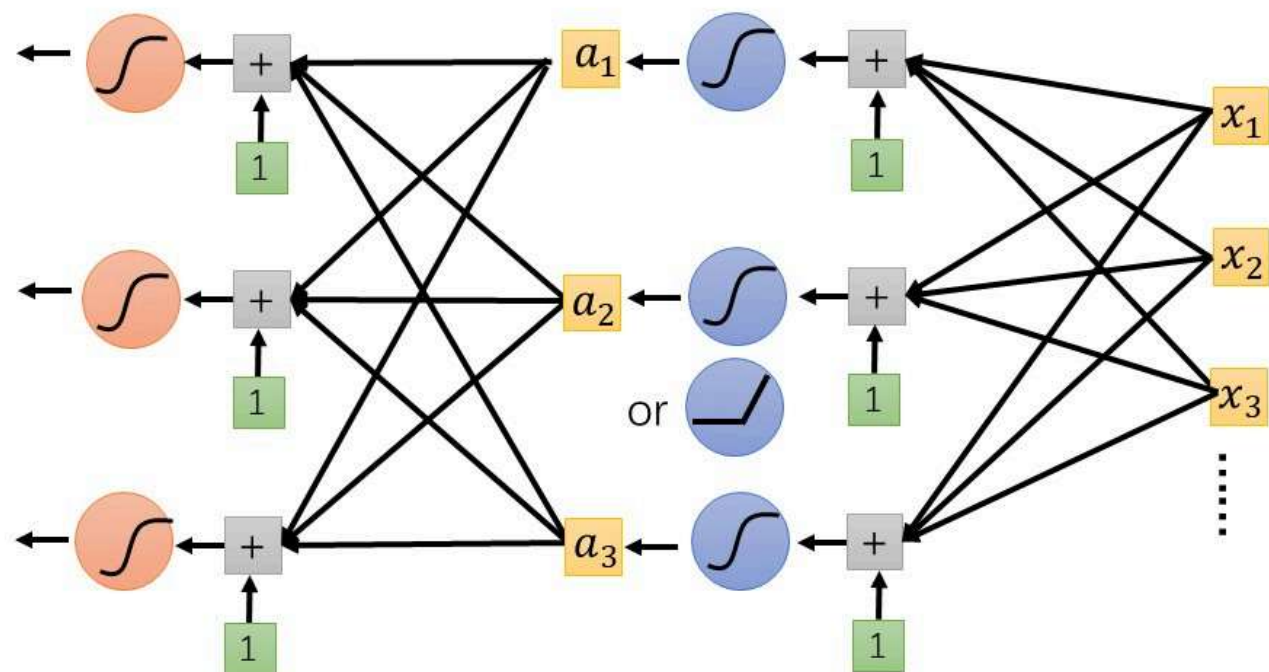
$$y = b + \sum_i c_i \max\left(0, b_i + \sum_j w_{ij} x_j\right)$$

	linear
2017 – 2020	0.32k
2021	0.46k

回到深度学习框架



Even more variety of models ...



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

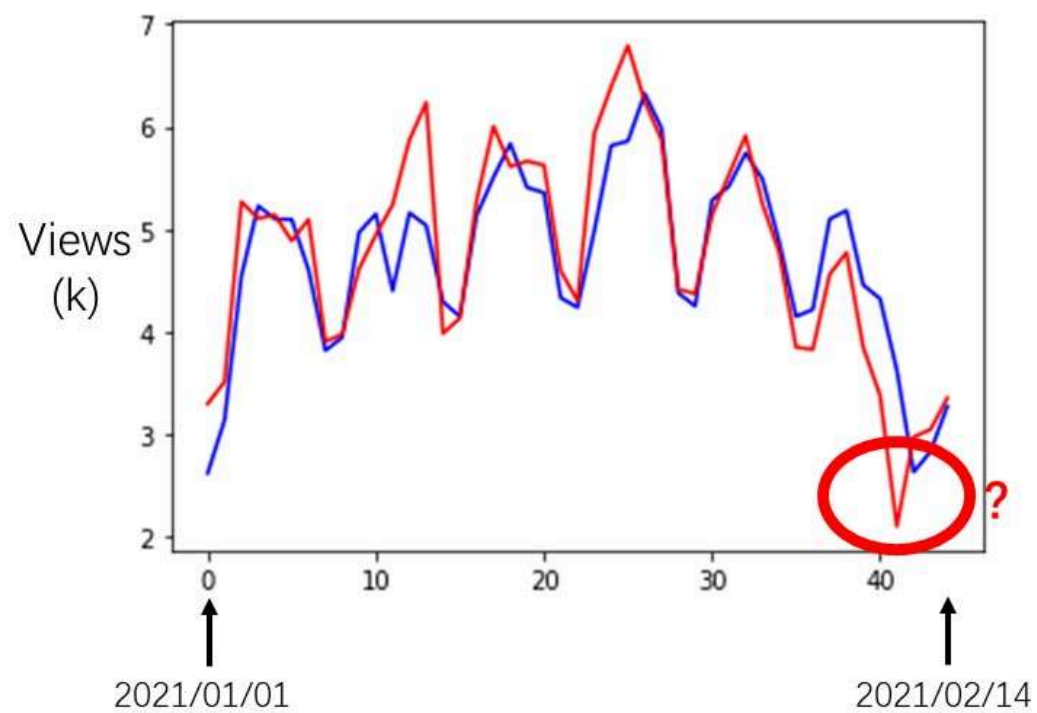
实验结果

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

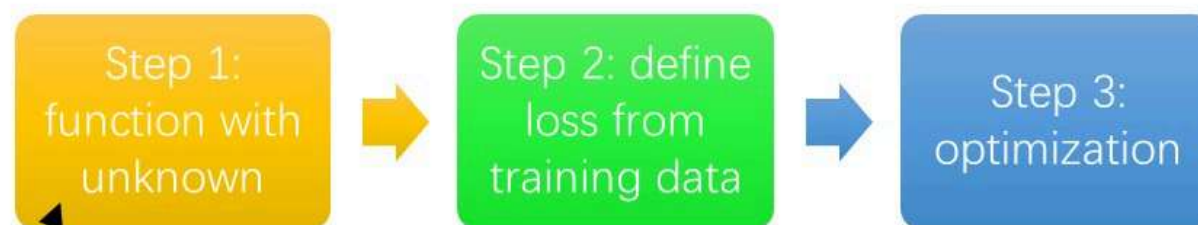
	1 layer				
2017 – 2020	0.28k				
2021	0.43k	0.35k			

3 layers

Red: 实际观看人数
blue: 预测观看人数



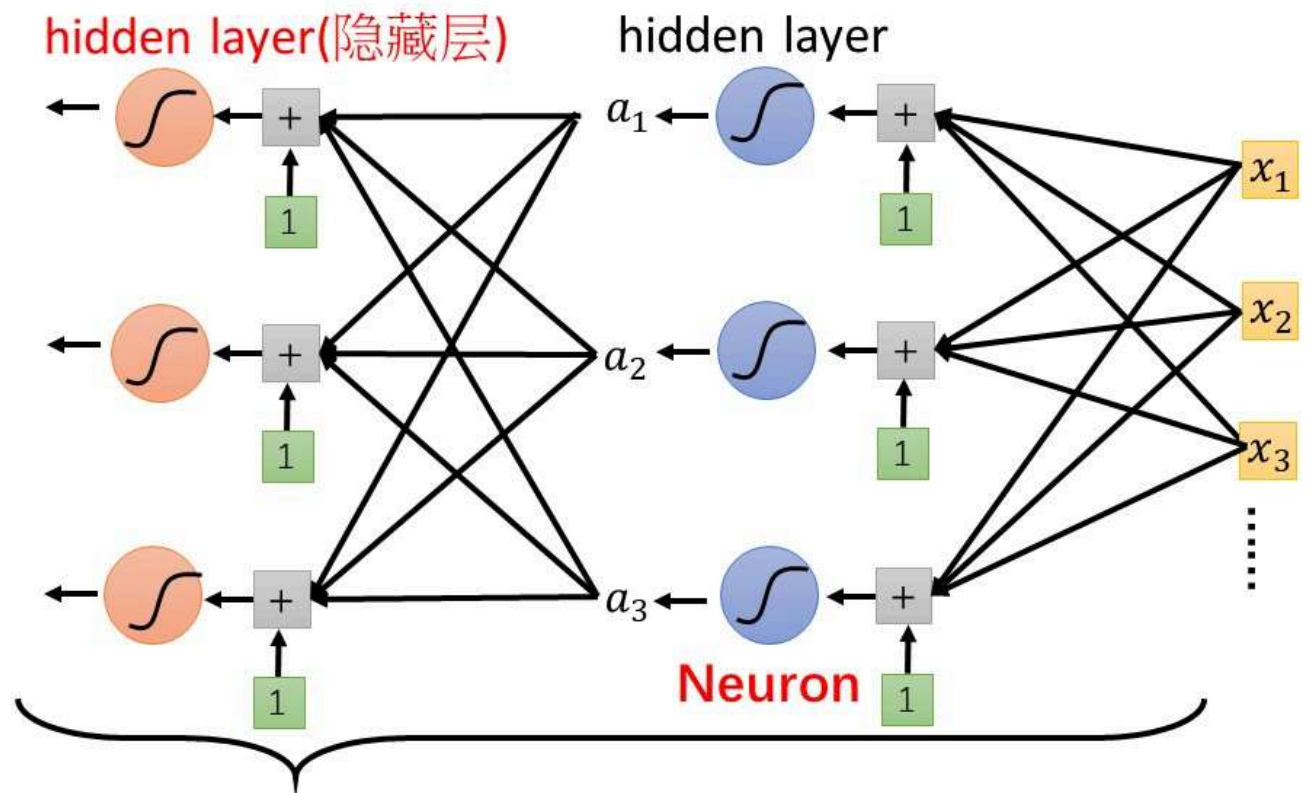
Back to DL Framework



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

It is not *fancy* enough.

Let's give it a *fancy* name!

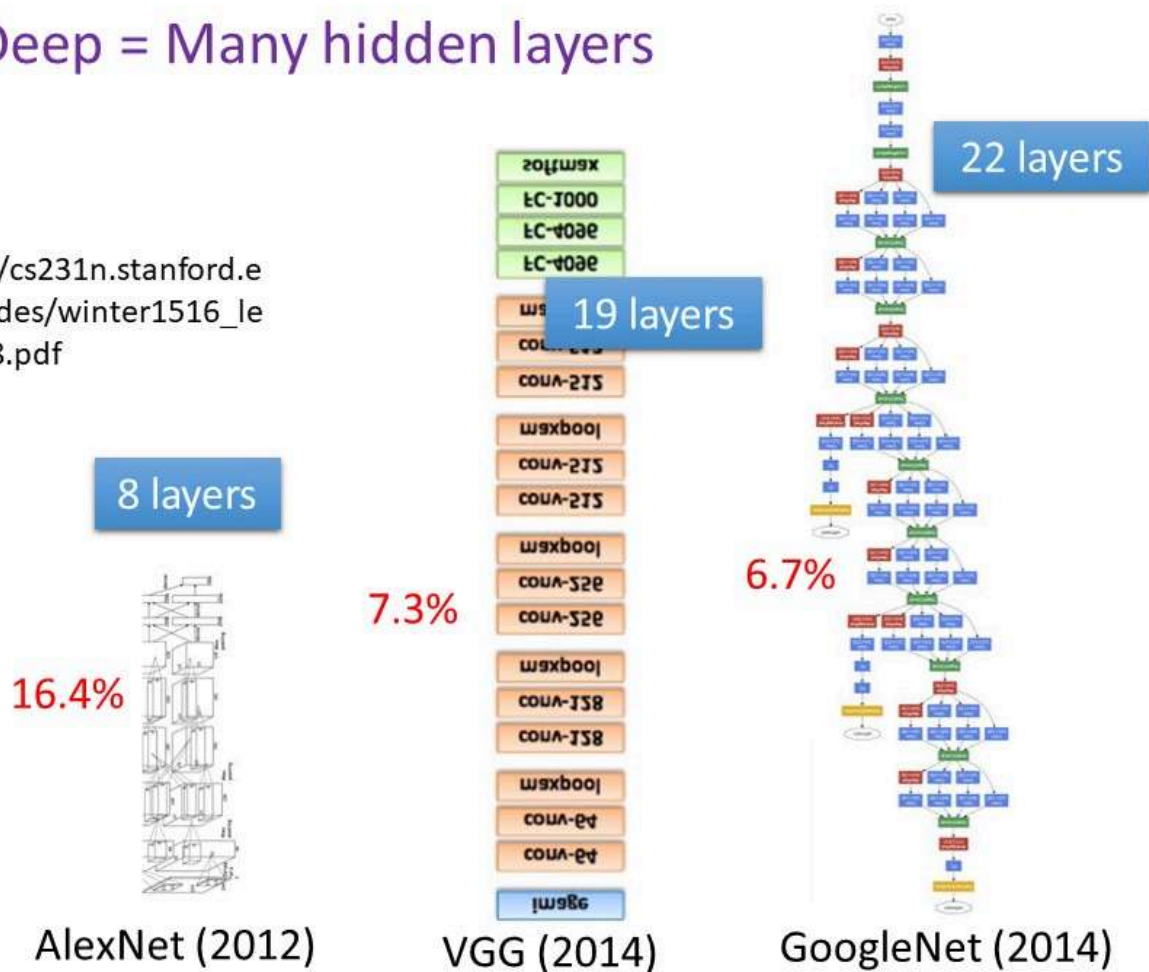


Neural Network This mimics human brains ... (???)

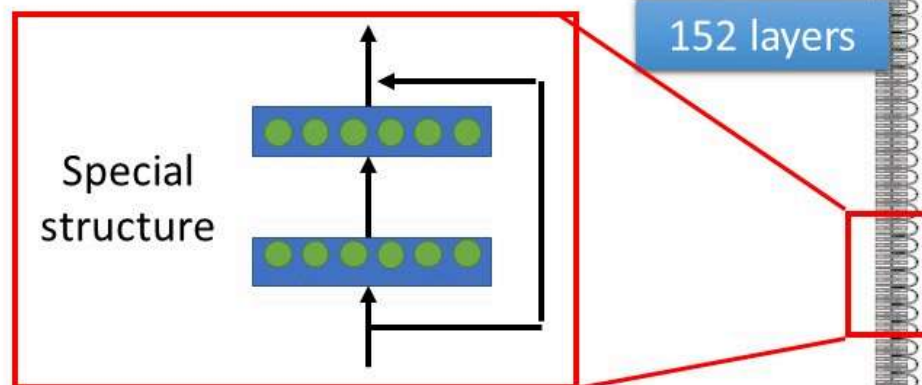
Many layers means **Deep** ➡ **Deep Learning**

Deep = Many hidden layers

http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf



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)

150 ~ layers



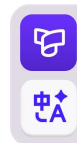
Guangzhou
Tower

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

0.44k



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**

Let's predict no. of views today!

- If we want to select a model for predicting no. of views today, which one will you use?

	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

We will talk about model selection next time. 😊