



机器学习和深度学习基本概念简介

机器学习

机器学习=寻找函数

不同的函数：regression（输出scale），classification（选择），Structured learning

如何寻找函数：

- ▼ 1、写出一个带有未知数的函数（base on domain knowledge）（也许是个猜测）

$y=b+wx_1$ ，这里带有未知数的函数为model， x_1 为Feature， w 为权重（weight）和 b 为偏置（bias）

- ▼ 2、定义Loss， $L(b, w)$ 。

用于评价 w 和 b 的设置是好还是不好，label是指真实值。

损失的计算方式有MAE（绝对差），MSE（相减平方）。

调整参数画出等高线图，是Error Surface。

- ▼ 3、最佳化。

找合适的未知参数的值使Loss最小。Gradient Descent（梯度下降）。

梯度下降：

假设只有一个参数 w ， w 不同数值得到不同Loss，获得曲线。

随机选取 w 的初始值 w_0 ，然后计算 w_0 的微分。 $\frac{\partial L}{\partial w} \Big|_{w=w_0}$

向Loss下降的方向移动，移动距离由 $\eta \frac{\partial L}{\partial w} \Big|_{w=w_0}$ 决定，这里的 η 是learning rate（学习速率），是自己设定的。自己决定的参数是hyperparameters。

移动结束有两个结果，一个是达到运动次数停止，理想停下是微分为0时。但有时候无法找到全局最优，大部分为局部最优（这是假命题）。

因此由一个参数 w 推广为两个参数 w, b ，步骤为：

1. 初始化：随机选取 w_0, b_0

2. 求微分偏导: $\frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0}, \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$
3. 更新w, b: $w^1 \leftarrow w^0 - \eta \frac{\partial L}{\partial w} \Big|_{w=w^0, b=b^0}, b^1 \leftarrow b^0 - \eta \frac{\partial L}{\partial b} \Big|_{w=w^0, b=b^0}$

以上步骤可以归结为训练

深度学习

Linear Model太过于简单, 无法模拟折线曲线, 也叫**Model bias**

所有的piecewise Linear curves (由多个锯齿状的线段组成的) 可

All Piecewise Linear Curves

= constant + sum of a set of



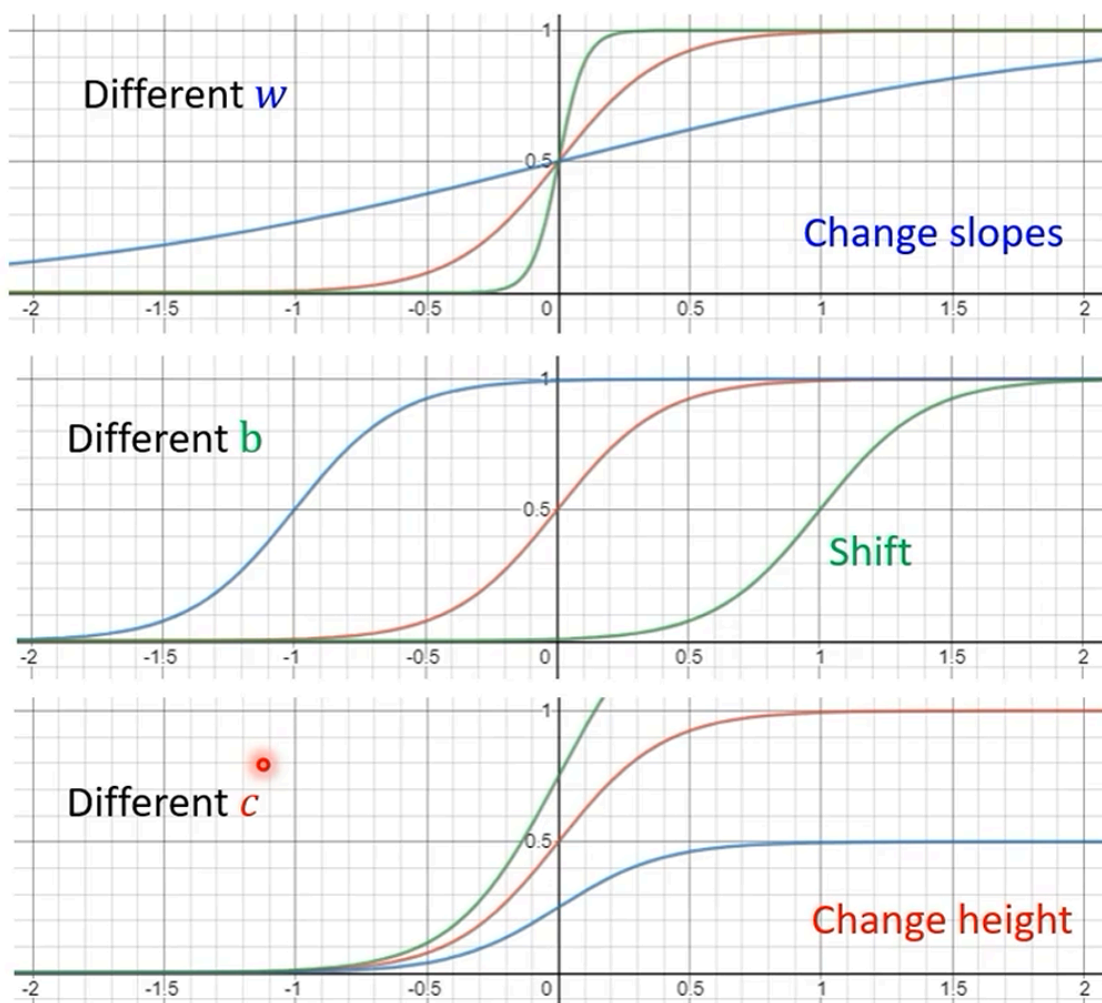
对于曲线, 可在曲线上取一些点再连接起来, 因此可以用piecewise Linear curves 去逼近任意一条曲线, 而piecewise Linear curves 可以用足够多的蓝色函数去表示。

▼ 如何写出蓝色function? →用sigmoid function 去逼近它

sigmoid函数: $y = c \cdot \frac{1}{1+e^{-(b+wx_1)}} = c \cdot \text{sigmoid}(b + wx_1)$

当 x_1 趋于无穷大, 则 y 趋于 c , x_1 趋于无穷小, y 趋于0

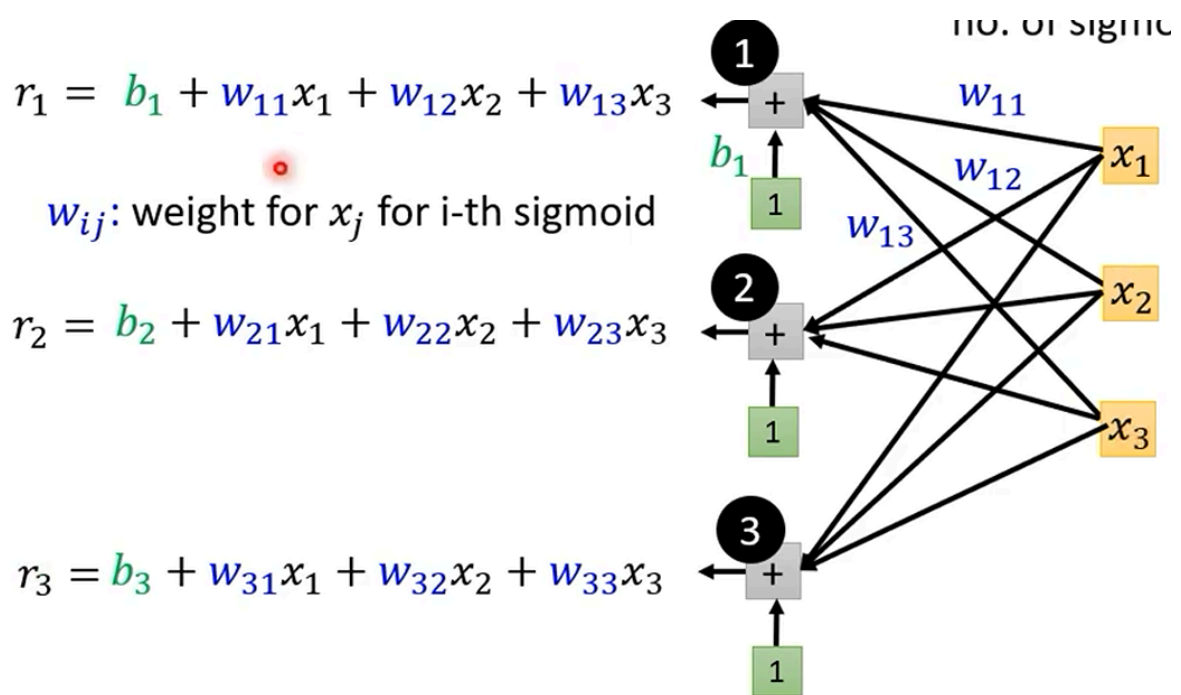
由此上面蓝色函数为hard sigmoid



因此，我们现在就有了如下的公式

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

这里的i是代表第i段sigmoid函数，j代表不同的Feature，就是例子中的考虑前三天的播放量的前三天数值。

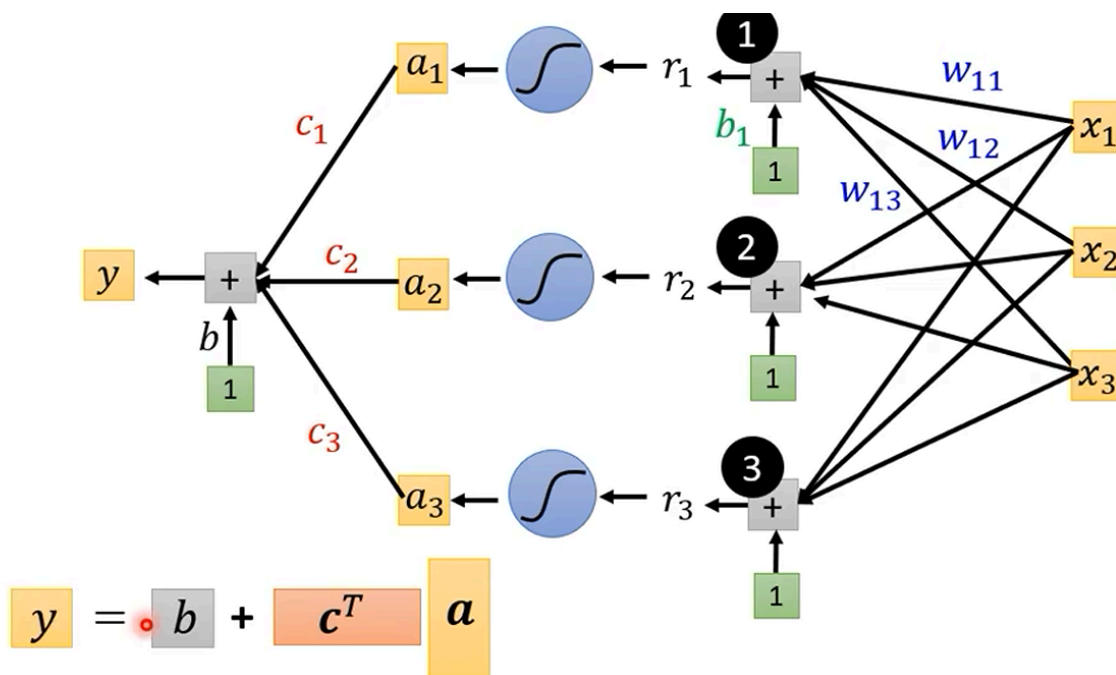


可以写为矩阵的形式：

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

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

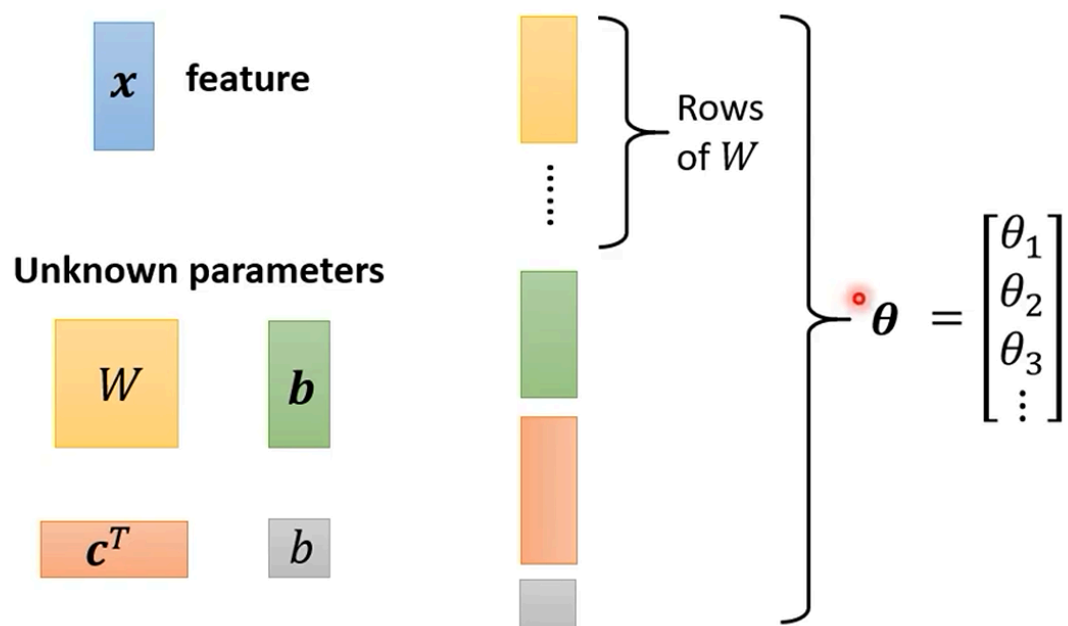
$$\mathbf{r} = \mathbf{b} + \mathbf{W} \mathbf{x}$$



其实下图中的W展开可以按row，也可以按column

Function with unknown parameters

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



由此上面就改写了机器学习的第一步，就是定义一个有未知数的函数
具体用几个sigmoid函数，是一个hyperparameters。

接下来定义Loss, $L(\theta)$, 跟之前没什么区别

然后是优化, 还是梯度下降

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

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

对于有大量数据data, 总数为N, 我们将其分为N/B个batch组, 每一个batch中有B个数据, 原本用所有数据计算Loss, 现在只选取一个B组中的数据进行计算

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

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

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

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

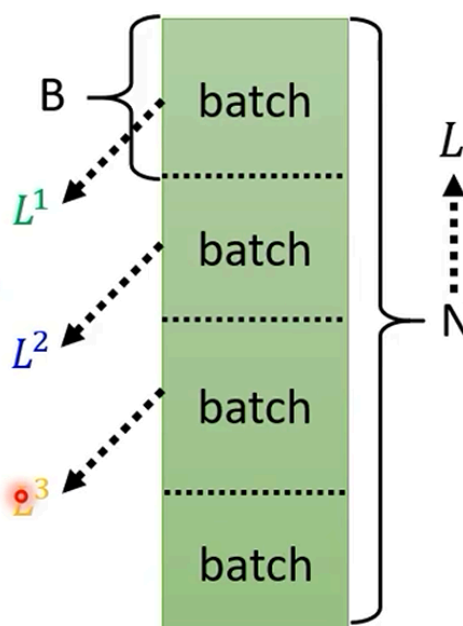
➤ Compute gradient $\mathbf{g} = \nabla L^2(\boldsymbol{\theta}^1)$

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

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

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

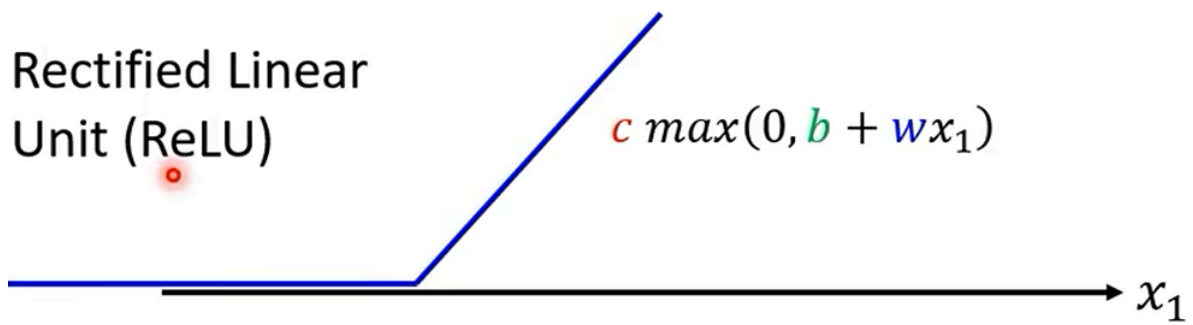
1 epoch = see all the batches once



每一次更新参数为update, 并且batch的大小也是自己决定的。

对于多个数据, 损失函数可以是所有的数据损失加起来, 也可以取平均

hard sigmoid可以由俩个rectified Linear Unit (ReLU) 组成



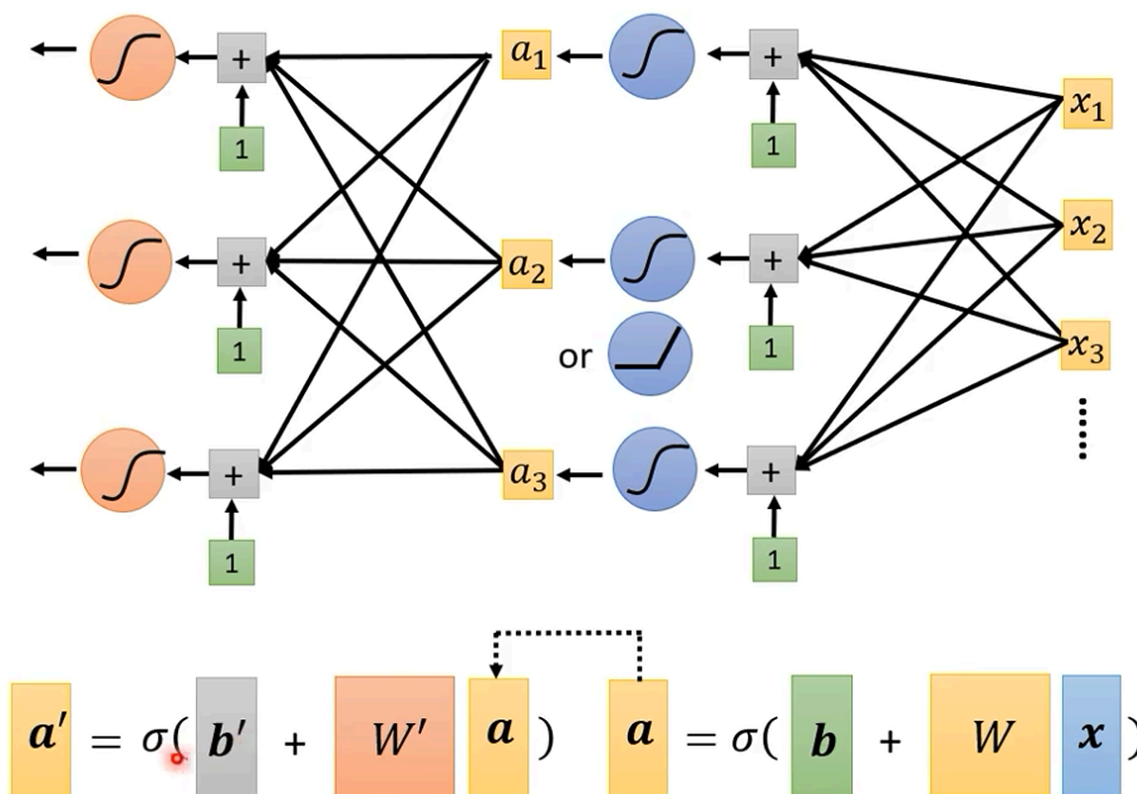
Sigmoid \rightarrow ReLU

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

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

注意由sigmoid变为ReLU需要将i变为2倍，这两个函数统称为activation function（激活函数）。

进一步改进模型可以采用多层：



sigmoid或ReLU叫Neuron，很多Neuron就是Neuron Network，后面将每一层叫为hidden layer，很多layers就是 Deep learning

对于训练数据和测试数据机器表现不一样的情况为overfitting，指在训练数据上变好，但是测试数据上变差