

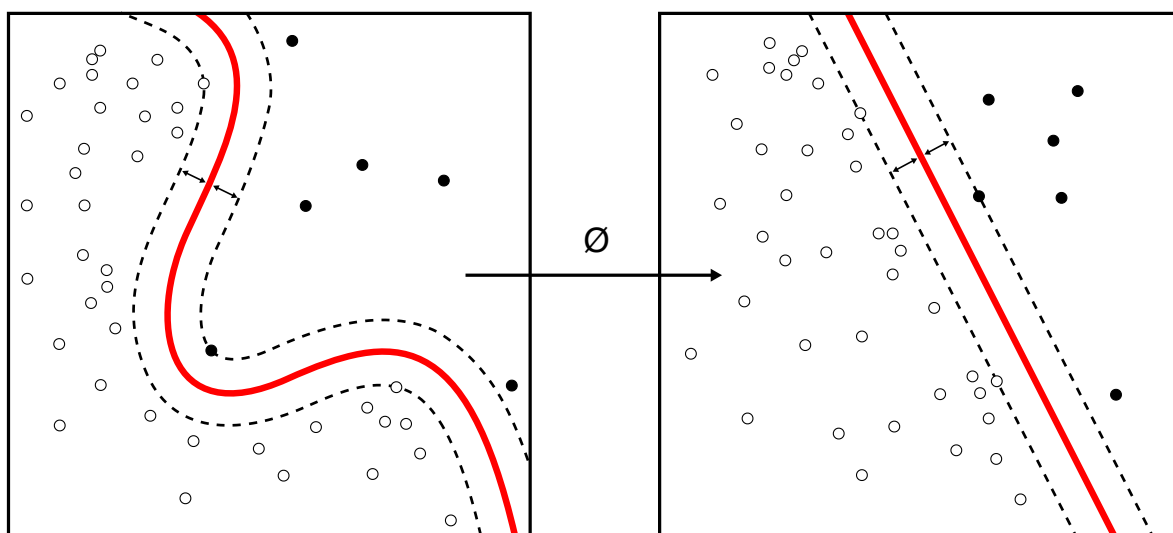
# Machine Learning

使用平台：作业Kaggle，免费GPU资源Google Colab

## 第一节：基本概念

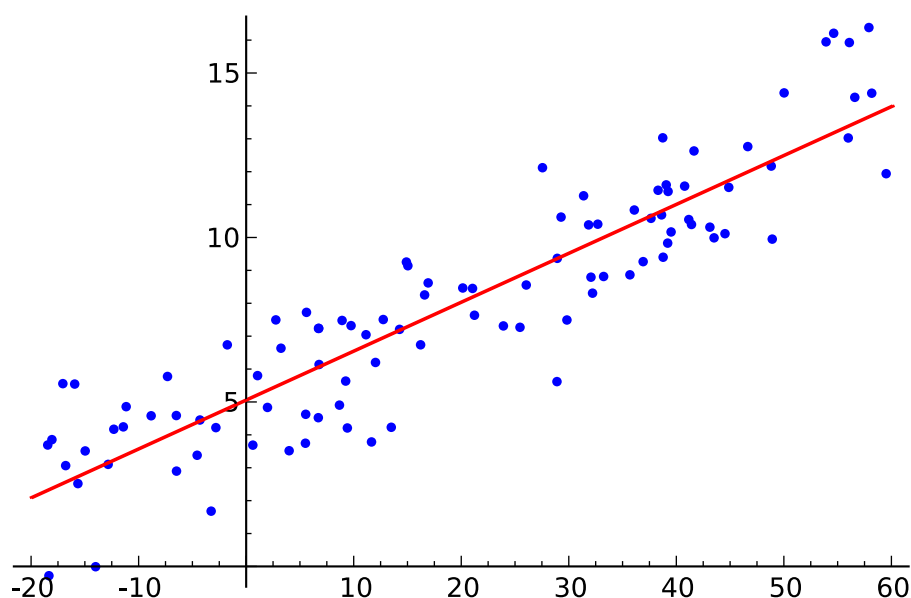
监督学习主要做两类工作，一是分类Classification，二是回归Regression

### 1. 分类Classification



我们试图将输入变量映射到离散类别中。例如给予患有肿瘤的患者，我们必须预测肿瘤是恶性的还是良性的。

### 2. 回归Regression

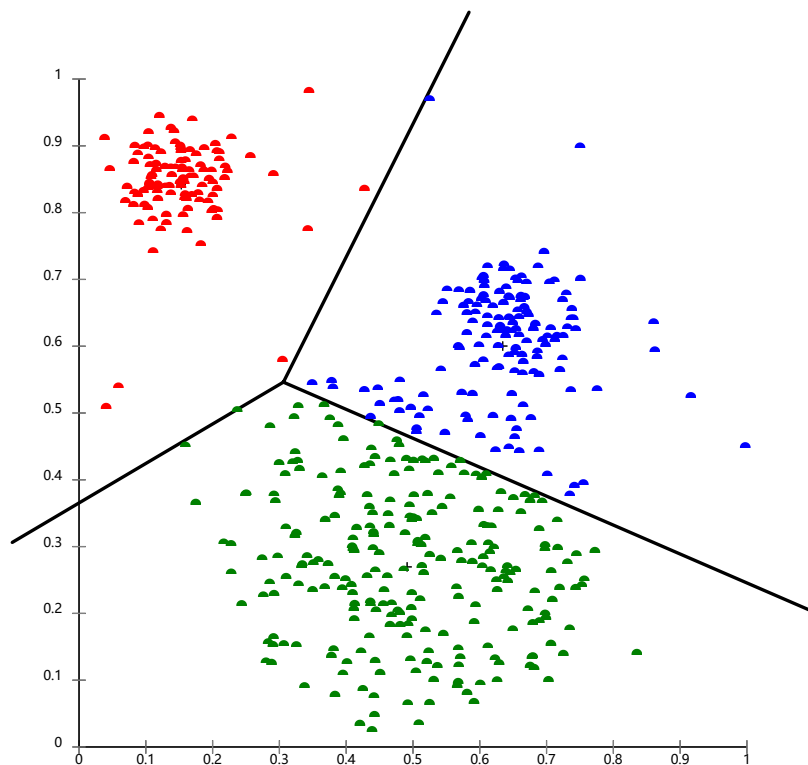


我们试图预测连续输出中的结果，这意味着我们试图将输入变量映射到某个连续函数。例如给定一个人的照片，根据照片预测年龄，这就是一个回归的问题。

非监督学习主要工作是“聚类”与“非聚类”

### 1. 聚类

获取1,000,000个不同基因的集合，并找到一种方法将这些基因自动分成不同变量的相似或相关的组，例如寿命，位置，角色等。



## 2. 非聚类

“鸡尾酒会算法”，允许您在混乱的环境中查找结果。（即在鸡尾酒会上识别来自声音网格的个人声音和音乐）

Machine Learning  $\approx$  Looking for function

（比如给定一个图片/音频/etc.，输出是图中物体/语音/etc.，我们要做的就是找到这个函数）

在“分类”与“回归”以外，机器学习还有“**Structured Learning**”，不只是做选择题

Create sth with structure(image,document)

猜测函数的流程（以预测某个channel今天的点击量为例）：

### 1. 设一含有未知参数的函数

假设 $y=b+wx_1$

y：今日点击量， $x_1$ ：昨日点击量(feature)

w(weight)和b(bias)是未知参数（从数据中得出）

没有乘以x的是bias，乘以x的是weight

### 2. 定义Loss函数

$L(b,w)$ 输入为b和w  $\uparrow$

假设 $b=0.5k, w=1$

$y=0.5k+x_1$

用预测的值与真实值差的绝对值作为评估效果

求出每一天的误差 $e_i$

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

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

### 3. Optimization(最优化)

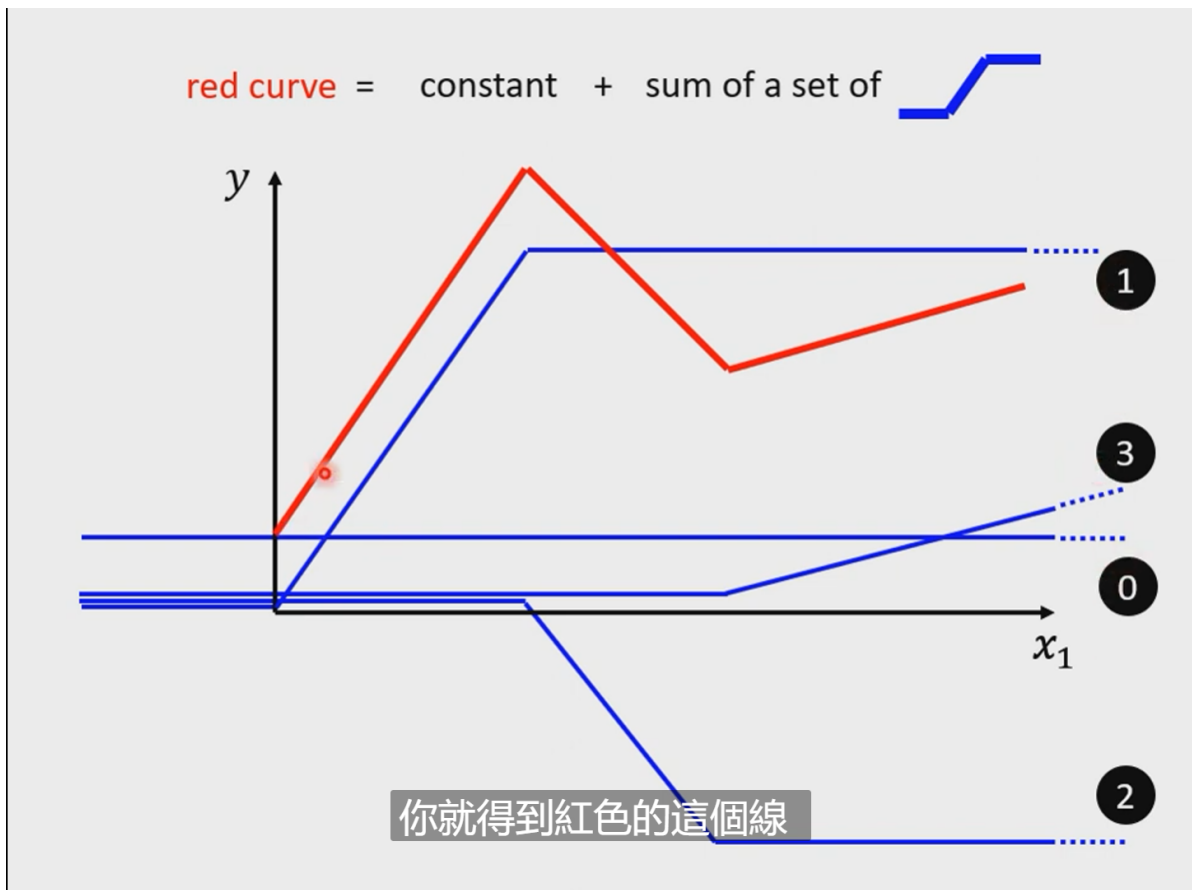
$$w^*, b^* = \arg \min(L)$$

#### Gradient Descent(梯度下降法)

- 随机选取初始值 $w^0$ 、 $b^0$
- 计算 $\frac{\delta L}{\delta w}|_{w=w^0}$   $\frac{\delta L}{\delta b}|_{w=w^0}$
- 设步长step为 $\eta \frac{\delta L}{\delta w}|_{w=w^0, b=b^0}$ 、 $\eta \frac{\delta L}{\delta b}|_{w=w^0, b=b^0}$  (斜率越大处步长越大)  $\eta$ : learning rate, 是自己设置的
- 类似 $\eta$ 这些自己设置的就是hyperparameters
- $w^1 = w^0 - \eta \frac{\delta L}{\delta w}|_{w=w^0, b=b^0}$ ,  $b^1 = b^0 - \eta \frac{\delta L}{\delta b}|_{w=w^0, b=b^0}$ , 更新 $w$ 和 $b$ , 不断迭代

会有local minima和global minima, gradient descent会有local minima的问题 (在实际操作中是个假问题? ) ? ? ? ? ? ? ? ?

Linear model有model本身的限制, 真实函数图像可能非常复杂, 我们可以根据以下这种使用linear叠加的方法计算真实函数图像

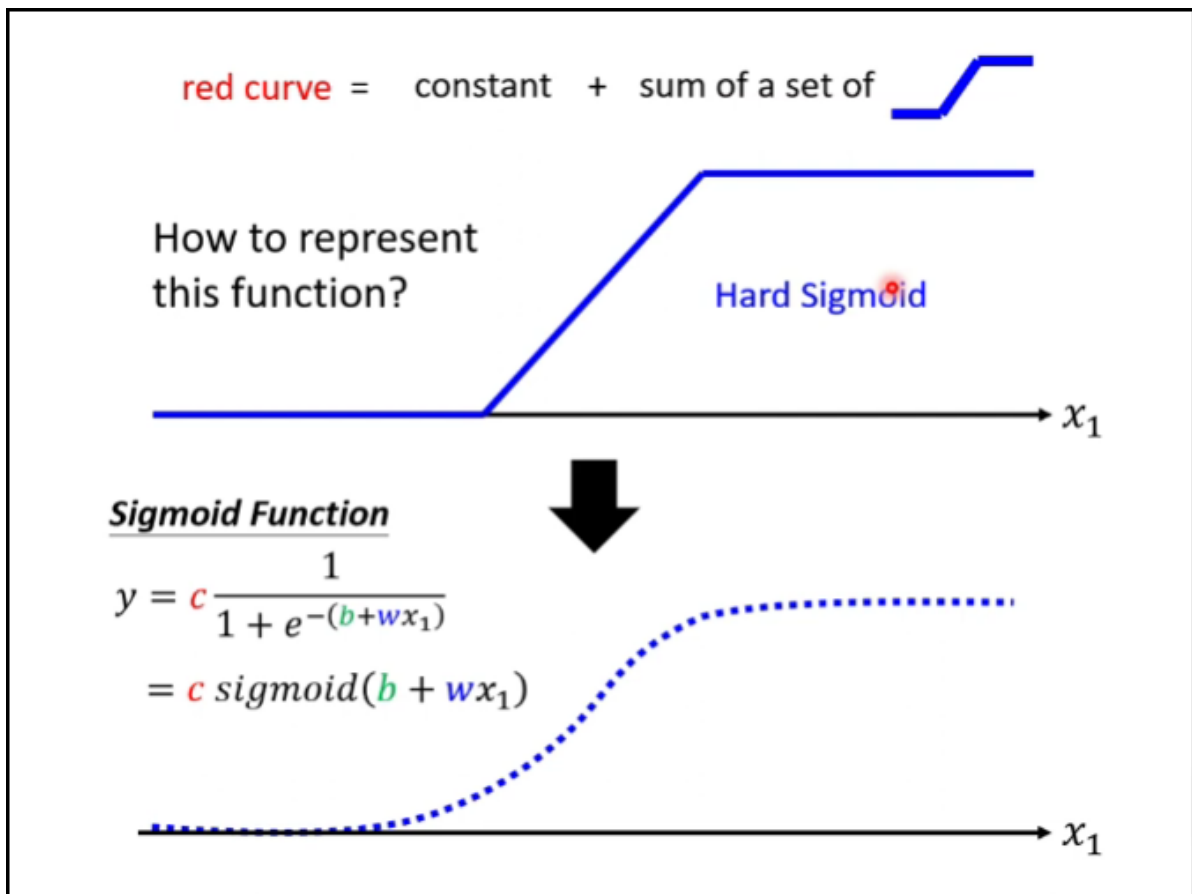


用这种方法可以表示任意函数

可以用Sigmoid函数逼近蓝色的函数

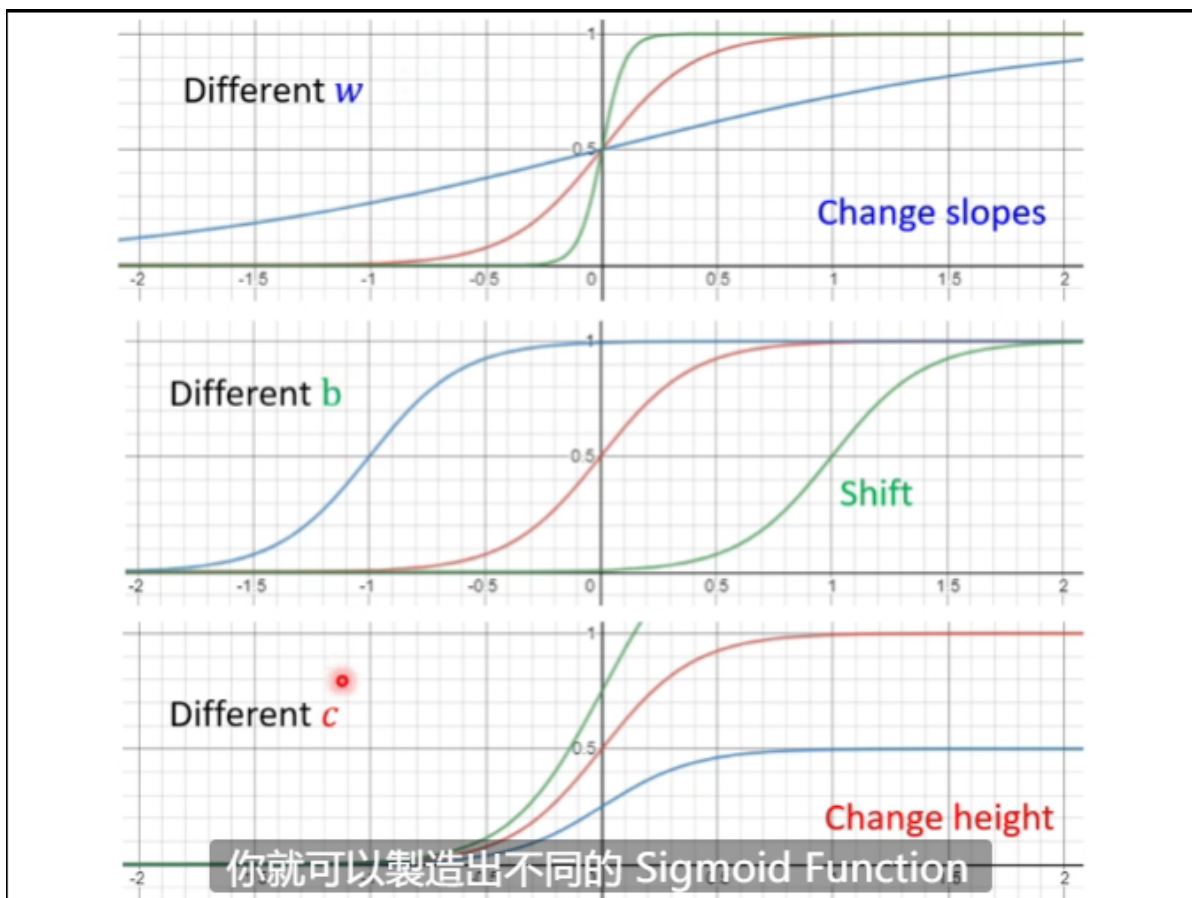
$$\text{Sigmoid函数为: } y = c \frac{1}{1 + e^{-(b + wx_1)}} = c * \text{sigmoid}(b + wx_1)$$

sigmoid可以翻译为S型的



蓝色函数叫做Hard Sigmoid

为了获得多种多样的函数，我们去调整b和w和c



More Features

## New Model: More Features

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

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

$$y = b + \sum_j w_j x_j$$

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

我們就把 Sigmoid 裡面的東西換掉

計算 $\theta$ 向量

## Optimization of New Model

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

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

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

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

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

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

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

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

是 0 向量 是 Zero Vector

## Optimization of New Model

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

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

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

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

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

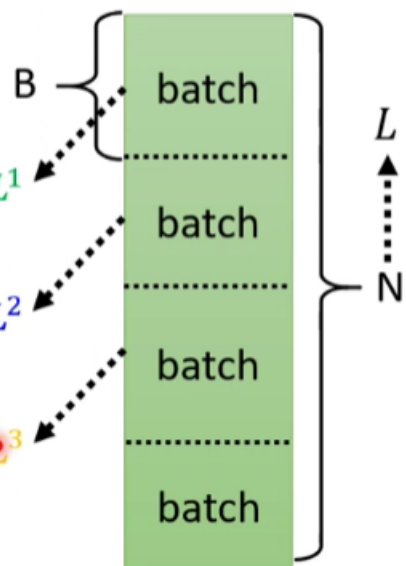
$$\text{update } \theta^2 \leftarrow \theta^1 - \eta g$$

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

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

1 epoch = see all the batches once

叫做一個 Epoch



## Optimization of New Model

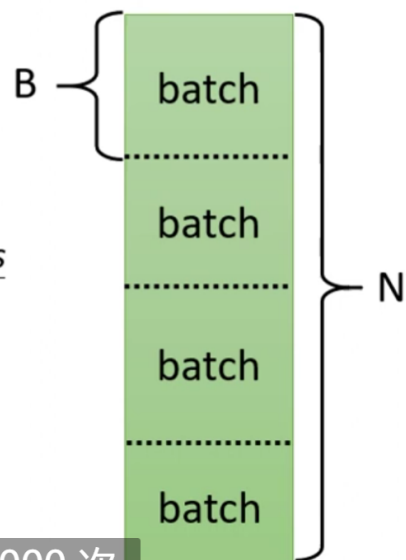
### Example 1

➤ 10,000 examples ( $N = 10,000$ )

➤ Batch size is 10 ( $B = 10$ )

How many update in 1 epoch?

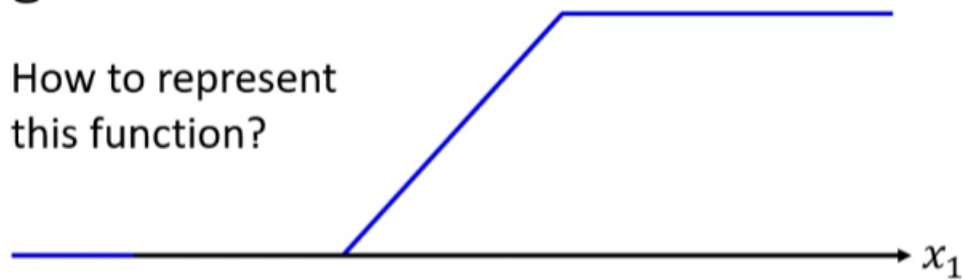
1,000 updates



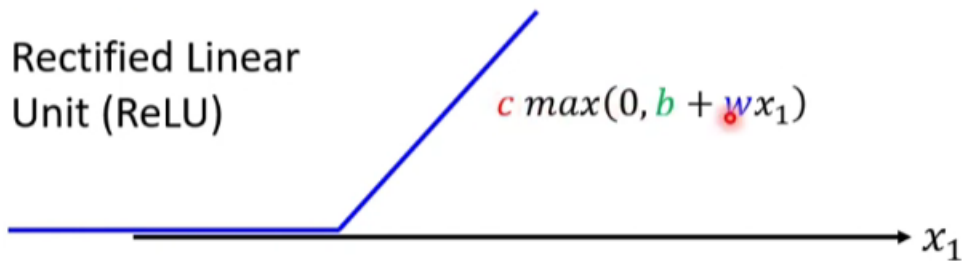
你其實已經更新了參數 1000 次

## Sigmoid → ReLU

How to represent  
this function?



Rectified Linear  
Unit (ReLU)



可以寫成  $c \max(0, b + wx_1)$

为什么ReLU比sigmoid好呢? ? ?

## Sigmoid → ReLU

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

Activation function

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

Which one is better?