



机器学习工程师直通车

—— 深度学习部分

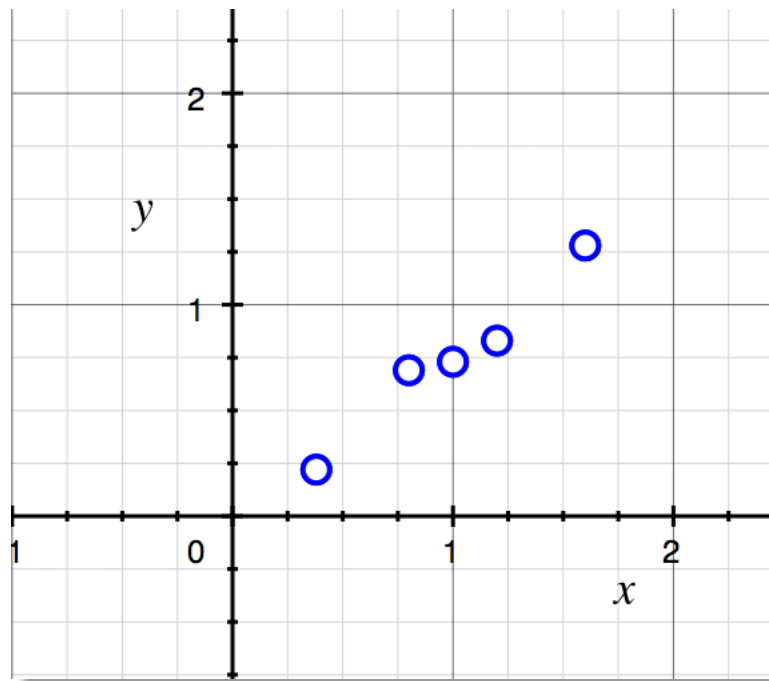
讲师：智亮

• 回顾一下梯度下降算法

监督学习

$$y = \theta x$$

x(input)	y(ground truth)
1	0.73
0.38	0.22
1.2	0.83
1.6	1.28
0.8	0.69



名词解释：

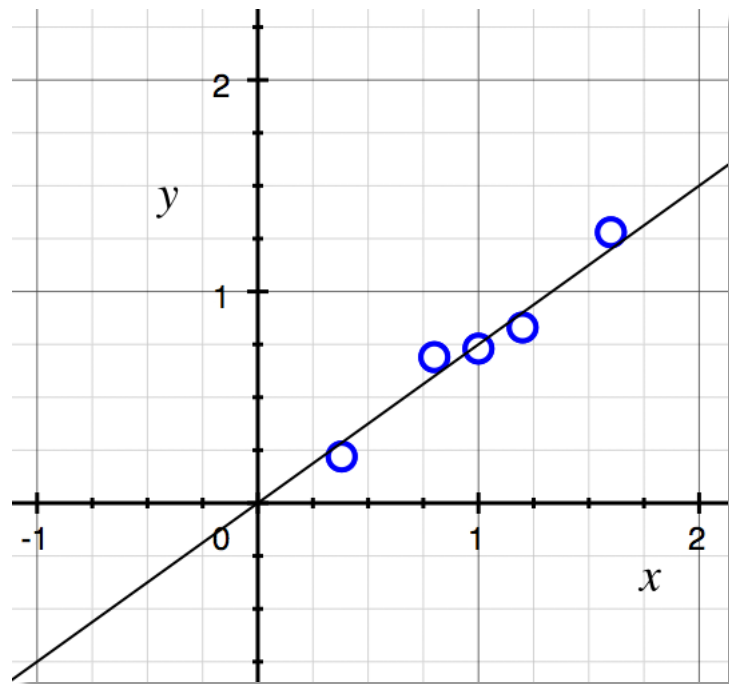
- 损失 (loss)：衡量神经网络的输出和Ground truth有多大差异的值。
- 损失函数 (loss function)：计算损失的函数。又称为代价函数，误差函数。
- 梯度 (gradient)：损失函数对于权重的导（函）数，用来衡量权重改变时，损失会如何变化。
- 梯度下降 (gradient descent)：根据梯度来更新权重，使损失变小的方法。
- 随机梯度下降 (stochastic gradient descent)：每次使用一条（或几条）数据（而不是全部数据），进行梯度下降的方法。
- 训练 (training)：使用梯度下降更新一点点 θ 。不断重复这个动作，使loss下降，直到获得最小值的过程。
- 训练step：一次“计算梯度，更新权重”的动作。
- 数据集 (dataset)：包含输入数据和ground truth，用来训练神经网络的数据。
- epoch：在随机梯度下降中，当数据集里面每条数据都用来训练了一次，叫一个epoch。

- 回顾一下梯度下降算法

监督学习

$$y = \theta x$$

x(input)	y(ground truth)
1	0.73
0.38	0.22
1.2	0.83
1.6	1.28
0.8	0.69



误差和损失函数 (loss/cost function)

$$\begin{aligned} loss &= \frac{1}{2}(\text{output} - y)^2 \\ &= \frac{1}{2}(\theta x - y)^2 \end{aligned}$$

- 损失 (loss)：衡量神经网络的输出和Ground truth有多大差异的值。
- 损失函数 (loss function)：计算损失的函数。又称为代价函数，误差函数。

- 回顾一下梯度下降算法

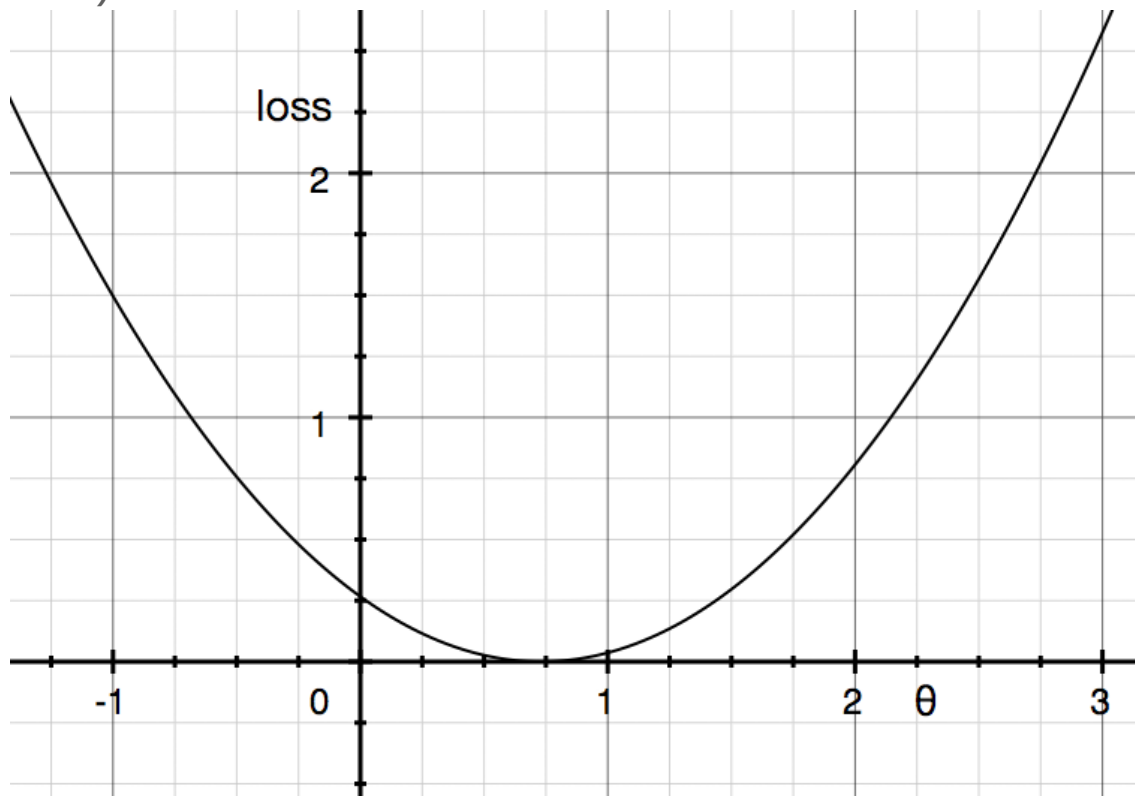
- 监督学习
- 误差和损失函数 (loss/cost function)

x(input)	y(ground truth)
1	0.73

$$y = \theta x \quad t = (\text{ground truth})$$

$$\begin{aligned} \text{loss} &= \frac{1}{2}(t - y)^2 \\ &= \frac{1}{2}t^2 - ty + \frac{1}{2}y^2 \\ &= \frac{1}{2}t^2 - t\theta x + \frac{1}{2}\theta^2 x^2 \end{aligned}$$

$$L(\theta) = \frac{1}{2}x^2 \cdot \theta^2 - tx \cdot \theta + \frac{1}{2}t^2$$



- 损失 (loss)：衡量神经网络的输出和Ground truth有多大差异的值。
- 损失函数 (loss function)：计算损失的函数。又称为代价函数，误差函数。

• 回顾一下梯度下降算法

- 监督学习
- 误差和损失函数 (loss/cost function)

$$y = \theta x \quad t = (\text{ground truth})$$

$$\begin{aligned} \text{loss} &= \frac{1}{2}(t - y)^2 \\ &= \frac{1}{2}t^2 - ty + \frac{1}{2}y^2 \\ &= \frac{1}{2}t^2 - t\theta x + \frac{1}{2}\theta^2 x^2 \end{aligned}$$

$$L(\theta) = \frac{1}{2}x^2 \cdot \theta^2 - tx \cdot \theta + \frac{1}{2}t^2$$

损失函数的导函数称为梯度函数，代表着损失函数在 θ 轴各点的斜率

$$\begin{aligned} \text{grad}_{\theta} &= \frac{dL}{d\theta} \\ &= \frac{dL}{d(t - y)} \frac{d(t - y)}{d\theta} \\ &= (t - y) \cdot (-x) \\ &= (y - t)x \end{aligned}$$

$$\begin{aligned} \text{grad}_{\theta} &= \frac{dL}{d\theta} \\ &= x^2\theta - tx \\ &= x(\theta x - t) \\ &= (y - t)x \end{aligned}$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \cdot \text{grad}_{\theta}$$

- 梯度 (gradient)：损失函数对于权重的导 (函) 数，用来衡量权重改变时，损失会如何变化。
- 梯度下降 (gradient descent)：根据梯度来更新权重，使损失变小的方法。

• 回顾一下梯度下降算法

- 监督学习
- 误差和损失函数 (loss/cost function)

$$y = \theta x \quad t = (\text{ground truth})$$

$$\begin{aligned} \text{loss} &= \frac{1}{2}(t - y)^2 \\ &= \frac{1}{2}t^2 - ty + \frac{1}{2}y^2 \\ &= \frac{1}{2}t^2 - t\theta x + \frac{1}{2}\theta^2 x^2 \end{aligned}$$

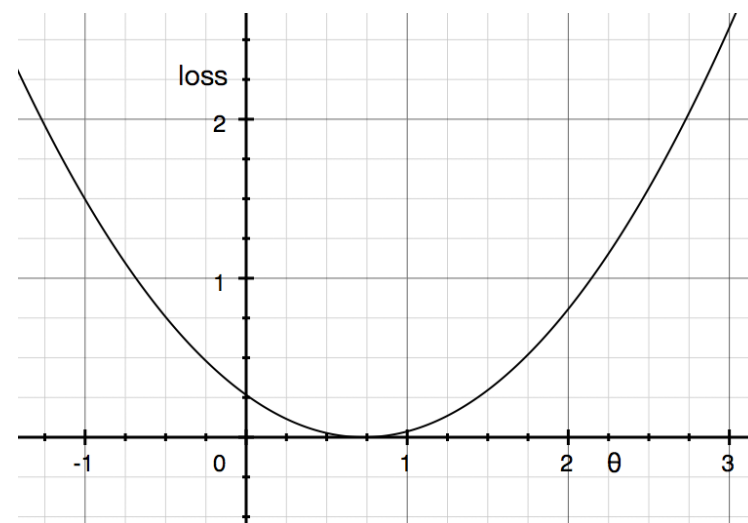
$$L(\theta) = \frac{1}{2}x^2 \cdot \theta^2 - tx \cdot \theta + \frac{1}{2}t^2$$

损失函数的导函数称为
梯度函数，代表着损失
函数在 θ 轴各点的斜率

$$\begin{aligned} \text{grad}_{\theta} &= \frac{dL}{d\theta} \\ &= \frac{dL}{d(t - y)} \frac{d(t - y)}{d\theta} \\ &= (t - y) \cdot (-x) \\ &= (y - t)x \\ \text{grad}_{\theta} &= \frac{dL}{d\theta} \\ &= x^2\theta - tx \\ &= x(\theta x - t) \\ &= (y - t)x \end{aligned}$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \cdot \text{grad}_{\theta}$$

x(input)	y(ground truth)
1	0.73



- 梯度 (gradient)：损失函数对于权重的导 (函) 数，用来衡量权重改变时，损失会如何变化。
- 梯度下降 (gradient descent)：根据梯度来更新权重，使损失变小的方法。

$$y = \theta x$$

$$loss = \frac{1}{2}(t - y)^2$$

$$grad_{\theta} = (y - t)x$$

$$\theta_{new} = \theta_{old} - \eta \cdot grad_{\theta}$$

若 $\eta=1$

1) $\theta = 3, x = 2, t = 8$

$$y = 3 \times 2 = 6$$

$$loss = \frac{1}{2}(8 - 6)^2 = 2$$

$$grad_{\theta} = (6 - 8) \times 2 = -4$$

$$\theta_{new} = 3 - (-4) = 7$$

2) $\theta = 7, x = 2, t = 8$

$$y = 7 \times 2 = 14$$

$$loss = \frac{1}{2}(8 - 14)^2 = 18$$

$$grad_{\theta} = (14 - 8) \times 2 = 12$$

$$\theta_{new} = 7 - 12 = -5$$

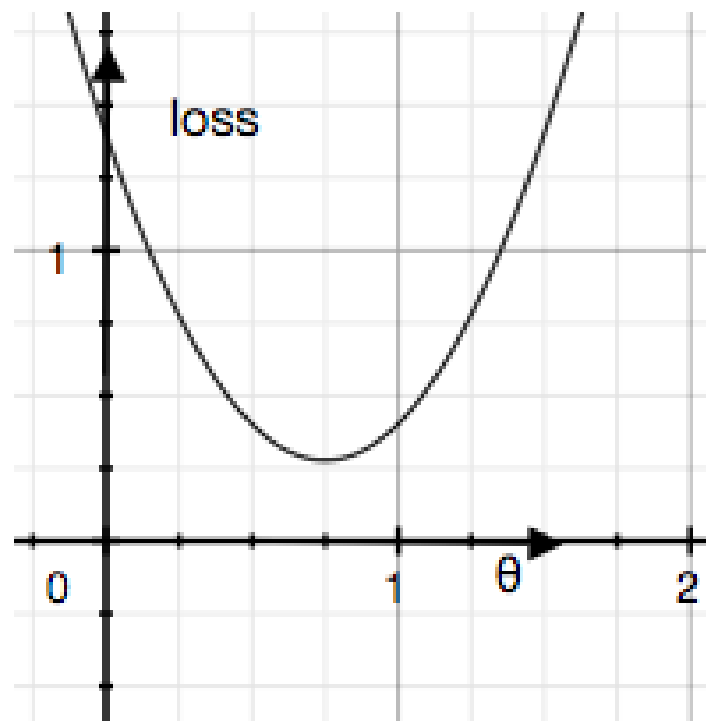
3) $\theta = -5, x = 2, t = 8$

$$y = -5 \times 2 = -10$$

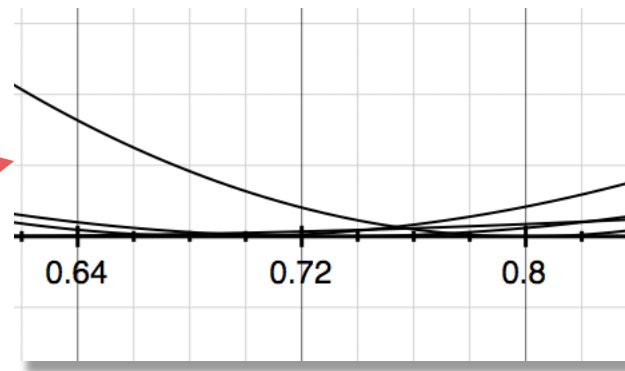
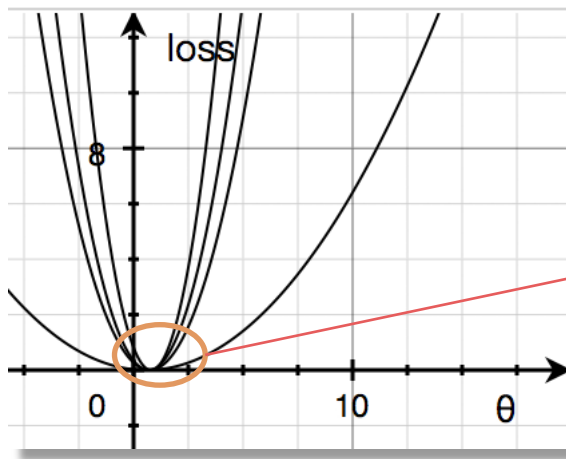
$$loss = \frac{1}{2}(8 - (-10))^2 = 162$$

$$grad_{\theta} = (-10 - 8) \times 2 = -36$$

$$\theta_{new} = -5 - (-36) = 31$$



x(input)	y(ground truth)
1	0.73
0.38	0.22
1.2	0.83
1.6	1.28
0.8	0.69

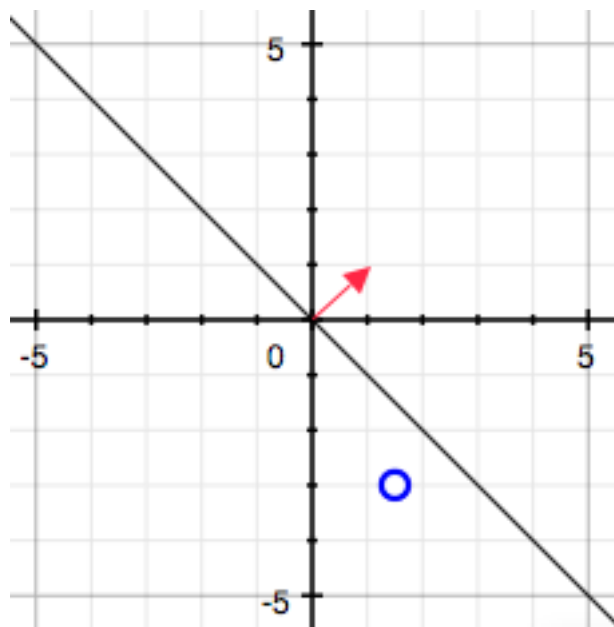


- 梯度下降 (gradient descent)：根据梯度来更新权重，使损失变小的方法。
- 随机梯度下降 (stochastic gradient descent)：每次使用一条（或几条）数据（而不是全部数据），进行梯度下降的方法。
- 训练 (training)：使用梯度下降更新一点点 θ 。不断重复这个动作，使loss下降，直到获得最小值的过程。
- 训练step：一次“计算梯度，更新权重”的动作。又称为iteration。
- 数据集 (dataset)：包含输入数据和ground truth，用来训练神经网络的数据。
- epoch：在随机梯度下降中，当数据集里面每条数据都用来训练了一次，称为一个epoch。

设 $x_1 = 1.5, x_2 = -3$
 $w_1 = 1, w_2 = 1$
 $b = 0$

此时 $output = f(w_1x_1 + w_2x_2 + b)$
 $= f(1 \times 1.5 + 1 \times -3 + 0)$
 $= f(-1.5)$
 $= 0$

约定 $y = output$, 若 $t=1$, 那么 $\Delta = t - y$
 $= 1 - 0$
 $= 1$



- 定义损失函数如下:

$$L(w, b) = -(t - y)(x \cdot w + b)$$

- 其对w和b的梯度为:

$$\text{grad}_w = \frac{dL}{dw} = -(t - y)x, \quad \text{grad}_b = \frac{dL}{db} = -(t - y)$$

- 对权重更新时使用:

$$w_{new} = w_{old} + \eta(t - y)x, \quad b_{new} = b_{old} + \eta(t - y)$$

- 对之前的例子应用权重更新:

$$\begin{aligned} W &= W + \eta(t - y)x \\ &= \begin{bmatrix} 1 & 1 \end{bmatrix} + 0.1 \times (1 - 0) \times \begin{bmatrix} 1.5 & -3 \end{bmatrix} \\ &= \begin{bmatrix} 1.15 & 0.7 \end{bmatrix} \end{aligned}$$

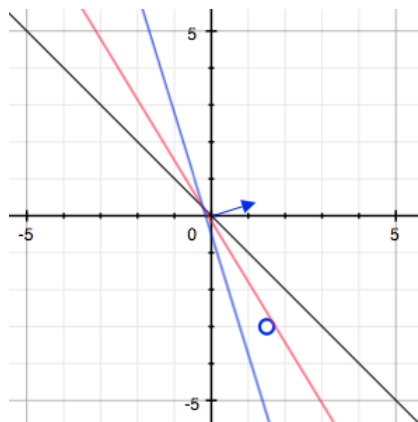
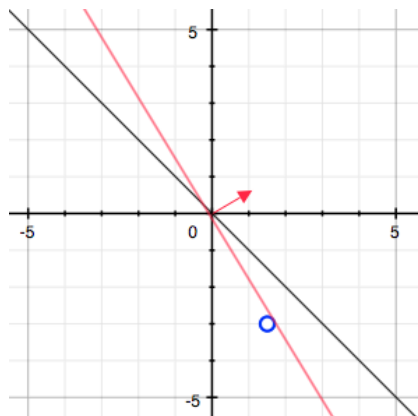
$$\begin{aligned} b &= b + \eta(t - y) \\ &= 0 + 0.1 \times (1 - 0) \\ &= 0.1 \end{aligned}$$

- 更新后的

$$\begin{aligned} \text{output} &= f(Wx + b) \\ &= f(1.15 \times 1.5 + 0.7 \times (-3) + 0.1) \\ &= f(-0.274) \\ &= 0 \end{aligned}$$

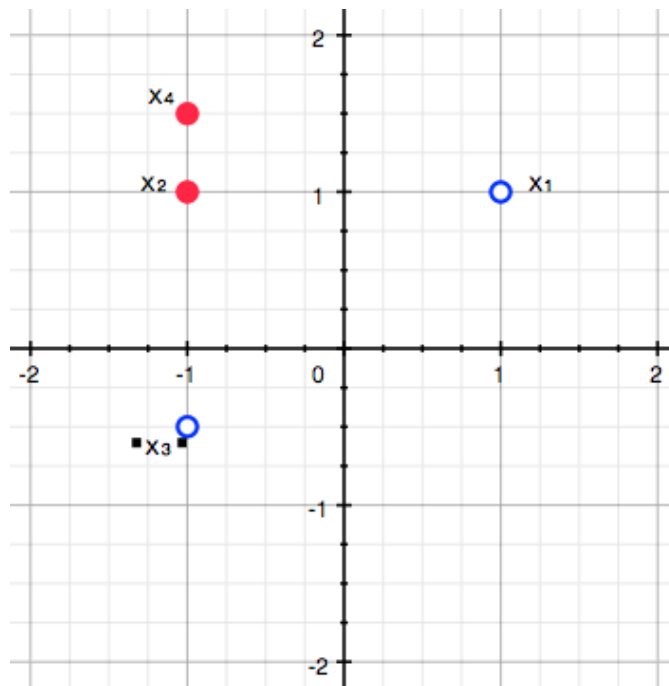
- 再次更新后

$$\begin{aligned} W &= \begin{bmatrix} 1.3 & 0.4 \end{bmatrix}, b = 0.2 \\ \text{output} &= f(1.3 \times 1.5 + 0.4 \times (-3) + 0.2) \\ &= f(0.95) \\ &= 1 \end{aligned}$$



- 多条数据情况
- 假定我们有这样几条输入：

$$X = \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ -1 & -0.5 \\ -1 & 1.5 \end{bmatrix} \quad t = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

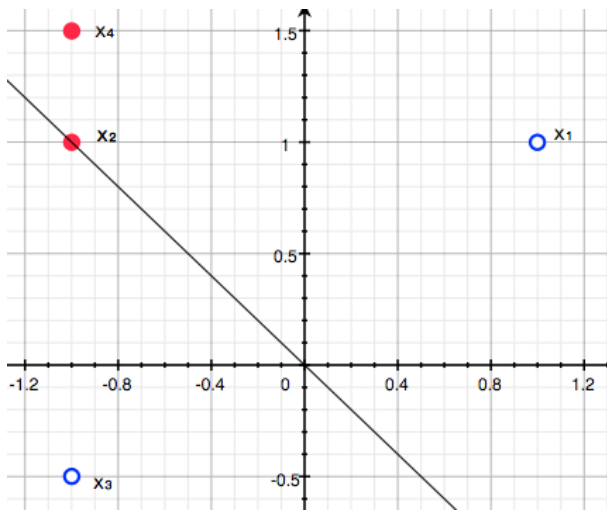


初始化设置 $w=[1, 1]$ ， $b=0$ ，则

$$\text{logits} = X \cdot w + b = \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ -1 & -0.5 \\ -1 & 1 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 0 = \begin{bmatrix} 1 \\ 0 \\ -1.5 \\ 0.5 \end{bmatrix}$$

$$\text{output} = f(\text{logits}) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \quad f(x) = \begin{cases} 1, x > 0 \\ 0, x \leq 0 \end{cases}$$

现在, $t = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix}, y = \text{output} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

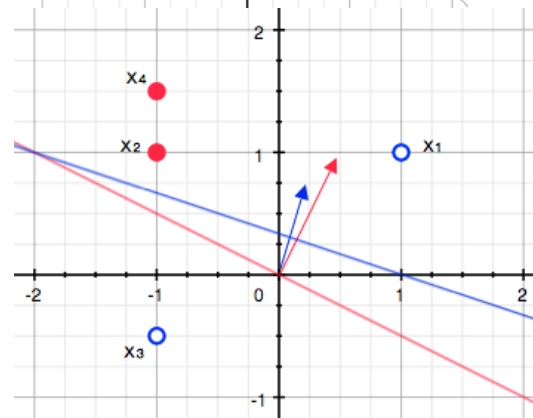
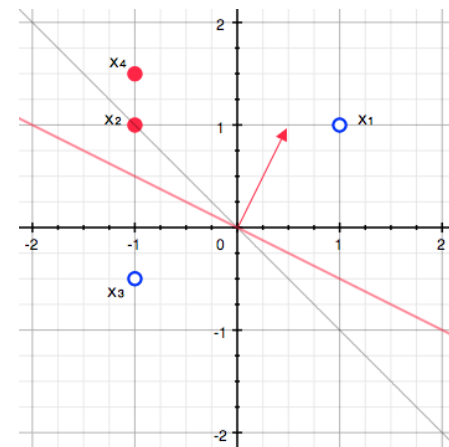


$$\begin{aligned}grad_w &= -\sum_i (t - y)x_i \\&= -\left(\begin{bmatrix} 0 \\ 1 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}\right) \cdot \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ -1 & -0.5 \\ -1 & 1 \end{bmatrix} \\&= -\begin{bmatrix} -1 \\ 1 \\ 0 \\ 0 \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ -1 & -0.5 \\ -1 & 1 \end{bmatrix} \\&= -[-2 \quad 0]\end{aligned}$$
$$\begin{aligned}grad_b &= -\sum (t - y) \\&= -(-1 + 1 + 0 + 0) \\&= 0\end{aligned}$$

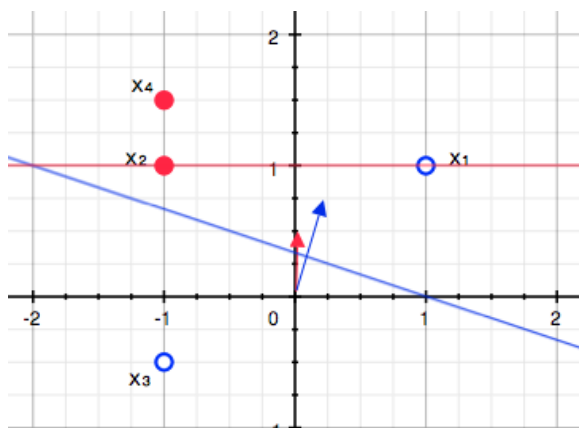
设 $\eta=0.25$, 则 $W_{new} = W_{old} - \eta grad_w = [0.5 \quad 1]$
 $b = 0$

重复这个过程, 得到新的

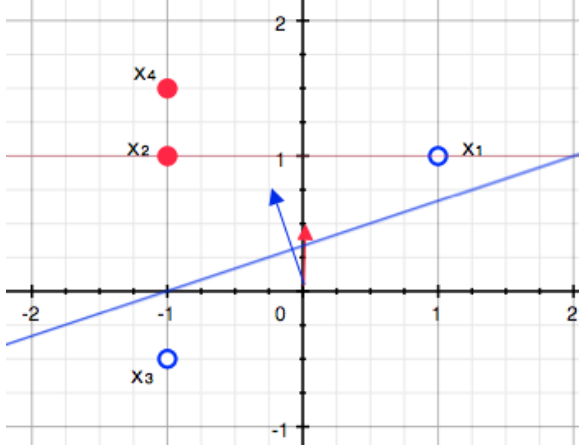
$$\begin{aligned}grad_w &= \sum_i -(t - y)x_i \\&= -[-1 \quad -1] \\W_{new} &= W_{old} - \eta grad_w = [0.25 \quad 0.75] \\b &= -0.25\end{aligned}$$
$$grad_b = \sum -(t - y) = 1$$



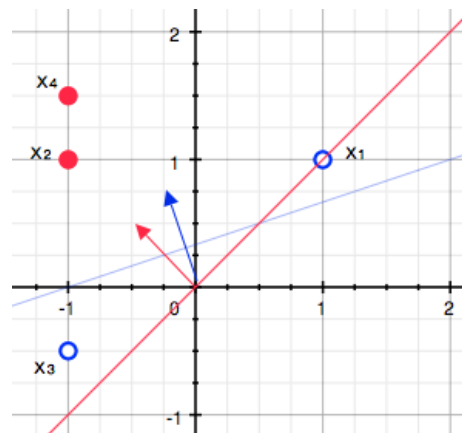
3) $W_{new} = W_{old} - \eta grad_w = [0 \quad 0.5]$
 $b = -0.5$



5)



4) $W_{new} = W_{old} - \eta grad_w = [0.25 \quad 0.75]$
 $b = -0.25$



6)

