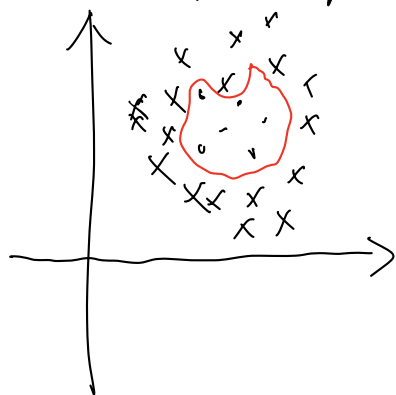


① Lecture 2 中的线性变化可以解决部分问题, 但有很大缺陷.

物理世界很多都是非线性问题, 比如 开或运算, 分等



红色曲线划分了两类数据, 但曲线本身是非线性

② 引入非线性激活函数  $h_o(w) = \theta^T \phi(x)$ ,  $\theta \in \mathbb{R}^{d \times k}$ ,  $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^d$ .

$\phi$  用于把  $x$  组成非线性的隐藏的特征

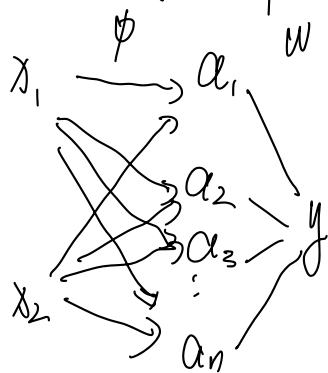
比如上面的曲线分类,  $\phi(x_1, x_2) = (x_1^2, x_2^2)$ ,  $\phi: \mathbb{R}^2 \rightarrow \mathbb{R}^2$

问题变成寻找合适的  $\phi$ , 映射出隐藏特征。这是古典机器学习:

人工识别数据特征, 设计  $\phi$ , 解决分类问题

但是, 这不具备通用性。场景之间  $\phi$  差异是巨大的

③ 无法提前选择适合的  $\phi$ , 那么只好全上!



$$\alpha_1 = x_1^2$$

$$\alpha_2 = x_2^2$$

$$\alpha_3 = \sin x_1$$

$$\alpha_4 = \cos x_1$$

$$\alpha_5 = \text{其它非线性函数}$$

预先准备常用的非线性函数, 由后

面的  $w$  来线性拟合, 通过学习  $w$  来

选择不同  $\phi$  的权重

④ 有了 feature mapping, 后面的层怎么办?

$$Z_1 \xrightarrow{W_1} Z_2 \xrightarrow{W_2} \dots \xrightarrow{W_L} Z_{L+1}$$

$$Z_{i+1} = G_i(Z_i; W_i), \rightarrow \begin{cases} i=1, 2, \dots, L \\ Z_i = X \\ h_0(X) = Z_{L+1} \end{cases}$$

$W$  即是网络参数

$G$  一般是 ReLU 激活函数

$$Z_i \in \mathbb{R}^{m \times n_i}, W_i \in \mathbb{R}^{n_i \times n_{i+1}}$$

⑤ 如何求  $W_i$ ?

$$\text{Loss} = \frac{1}{m} \sum_{i=1}^m \text{loss}(h_0(x^i), y^i)$$

$$\frac{d \text{Loss}}{dW}, \text{ 又成为求偏导的问题}$$

~~基于  $G$  是 softmax, 最后可以推导出~~  $\frac{d \text{Loss}}{dW_i} = \frac{d \text{Loss}}{dZ_{L+1}} \cdot \frac{dZ_{L+1}}{dZ_L} \dots \frac{dZ_{i+1}}{dW_i}$

此时我们看到了 **反向传播**,

记  $C_{i+1} = \frac{d \text{Loss}}{dZ_{L+1}} \cdot \frac{dZ_{L+1}}{dZ_L} \dots \frac{dZ_{i+2}}{dZ_{i+1}}$ , 那么  $\frac{d \text{Loss}}{dW_i} = C_{i+1} \cdot \frac{dZ_{i+1}}{dW_i}$

$C_{L+1} = \frac{d \text{Loss}}{dZ_{L+1}}$ , 而  $Z_{L+1} \rightarrow \text{softmax} \rightarrow \text{loss}$

$\therefore C_{L+1} = S - I_y$

$C_i = C_{i+1} \cdot \frac{dZ_{i+1}}{dZ_i} = C_{i+1} \cdot \frac{d G(Z_i; W_i)}{dZ_i} = C_{i+1} \cdot G'(Z_i; W_i) \cdot W_i$

$\frac{d \text{Loss}}{dW_i} = C_{i+1} \cdot \frac{dZ_{i+1}}{dW_i}$ , 而  $C_{i+1}$  由  $C_{i+2}$  给出, 即是反向传播!

⑥ 虽然有了反向传播的理念，我们还是需要基于一个easy例子来理解它

$$\begin{cases} h(x) = \sigma(XW_1)W_2 \end{cases}$$

$$\text{loss} = \frac{1}{m} \sum_{i=1}^m (\text{loss}(h(x_i); x, y_i))$$

$$\begin{array}{ccccccc} X & \xrightarrow{W_1} & \sigma & \xrightarrow{W_2} & H_2 & \xrightarrow{\text{softmax}} & Z \\ m \times n & n \times L & L \times 1 & L \times k & m \times k & & m \times k \end{array}$$

$$\text{Loss} = \frac{1}{m} \sum_{i=1}^m \log \sum_{i=1}^m \log z_{ij} y_i, \quad z_{ij} y_i = \frac{e^{h_{ij} y_i}}{e^{h_{i1}} + e^{h_{i2}} + \dots + e^{h_{ik}}}$$

$$\textcircled{7} \frac{d\text{loss}}{dW_2} = \frac{d\text{loss}}{dH_2} \cdot \frac{dH_2}{dW_2}$$

$$\begin{aligned} 1) \frac{d\text{loss}}{dH_{ij}} &= \frac{\frac{1}{m} d(\sum_{i=1}^m \log \frac{e^{h_{ij} y_i}}{e^{h_{i1}} + e^{h_{i2}} + \dots + e^{h_{ik}}})}{dh_{ij}} = \frac{1}{m} \cdot \frac{d \log \frac{e^{h_{ij} y_i}}{e^{h_{i1}} + e^{h_{i2}} + \dots + e^{h_{ik}}}}{dh_{ij}} \\ &= \frac{1}{m} \frac{d(\log e^{h_{ij} y_i} - \log(e^{h_{i1}} + e^{h_{i2}} + \dots + e^{h_{ik}}))}{dh_{ij}} \\ &= \frac{1}{m} (1 - \frac{e^{h_{ij} y_i}}{e^{h_{i1}} + e^{h_{i2}} + \dots + e^{h_{ik}}}), j = y_i \\ &\quad \frac{1}{m} (0 - \frac{e^{h_{ij} y_i}}{e^{h_{i1}} + e^{h_{i2}} + \dots + e^{h_{ik}}}), j \neq y_i \end{aligned}$$

$$\text{即是 } \frac{d\text{loss}}{dH_{ij}} = \frac{1}{m} (I - Z_{ij})$$

$$\text{推广到矩阵 } \frac{d\text{loss}}{dH_2} = \frac{1}{m} (I - Z)$$

$$2) \frac{dH_{ij}}{dW_2} = H_{ij} = H_{i1}W_{1j} + H_{i2}W_{2j} + \dots + H_{iL}W_{Lj}$$

$$H_{ij} \text{ 只与 } W_{aj} \text{ 有关, } \frac{dH_{ij}}{dW_{aj}} = H_{i1}, \frac{dH_{ij}}{dW_{2j}} = H_{i2}, \dots, \frac{dH_{ij}}{dW_{Lj}} = H_{iL}$$

$$\frac{dH_{ij}}{dW_{aj}} = H'_{ia}, \quad a \in \{1, 2, \dots, L\}, \quad \text{即是 } \frac{dH_2}{dW_2} = H_1$$

那么  $\frac{d\text{Loss}}{dW_2} = \frac{1}{m} (I - Z) \cdot H_1$ , 整理为  $\frac{d\text{Loss}}{dW_2} = \frac{1}{m} \cdot H_1^T (I - Z)$

$\begin{matrix} L \times L & & m \times k & & m \times L \end{matrix}$

⑧ 继续求  $\frac{d\text{Loss}}{dW_1}$   $\frac{d\text{Loss}}{dW_1} = \frac{d\text{Loss}}{dH_2} \cdot \frac{dH_2}{dH_1} \cdot \frac{dH_1}{dW_1}$

1)  $\frac{d\text{Loss}}{dH_2} = \frac{1}{m} (I - Z)$ , ⑦已有证明

2)  $\frac{dH_2}{dH_1} = W_2$  ( $H_2 = H_1 W_2$ )

$\begin{matrix} m \times k & m \times L & L \times k \end{matrix}$

3)  $\frac{dH_1}{dW_1} = \frac{d \sigma(XW_1)}{dW_1} = \frac{d \sigma(XW_1)}{d(XW_1)} \cdot \frac{d(XW_1)}{dW_1} = \sigma'(XW_1) \cdot X$

那么  $\frac{d\text{Loss}}{dW_1} = \frac{1}{m} (I - Z) \cdot W_2 \cdot \sigma'(XW_1) \cdot X$

$\begin{matrix} n \times L & & m \times k & & L \times k & & m \times L & & m \times n \end{matrix}$

整理得  $\frac{d\text{Loss}}{dW_1} = \frac{1}{m} X^T \left( (I - Z) \times W_2^T \circ \sigma'(XW_1) \right)$

$\begin{matrix} n \times L & & n \times m & & m \times k & & k \times L & & m \times L \end{matrix}$

矩阵乘积      标量积

⑨ 那么  $W_1/W_2$  有了 gradient descent

$$\left\{ \begin{aligned} \frac{d\text{Loss}}{dW_1} &= \frac{1}{m} X^T \left( (I - Z) W_2^T \circ \sigma'(XW_1) \right) \\ \frac{d\text{Loss}}{dW_2} &= \frac{1}{m} H_1^T (I - Z) \end{aligned} \right.$$