

神经网络基础

1. 二分类问题

· 样本 (x, y) , 其中 $x \in R^{n_x}$, $y \in \{0, 1\}$

· 样本数量为 m . $\{x^{(1)}, y^{(1)}\}, \{x^{(2)}, y^{(2)}\}, \dots, \{x^{(m)}, y^{(m)}\}$

· 数据的矩阵情况

$$X = \underbrace{\begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{bmatrix}}_m \quad n^x \quad y = \begin{bmatrix} y^{(1)}, y^{(2)}, \dots, y^{(m)} \end{bmatrix} \quad (R^{1 \times m})$$

2. Logistic Regression

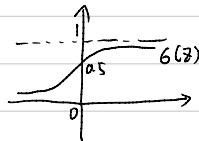
\hat{y} 表示预测的结果, $\hat{y} = P(y=1|x)$

输入: $X \in R^{n_x}$, $w \in R^{n_x}$, $b \in R$ 输出: $\hat{y} = w^T x + b$

引入 Sigmoid 函数 $\hat{y} = \text{Sigmoid}(w^T x + b)$, $\text{Sigmoid}(z) = \frac{1}{1+e^{-z}}$

存在偏置 b , 所以设置 $x_0 = 1$, 即 $x \in R^{n_x+1}$

参数 θ 也修改为 $\theta = [\underbrace{\theta_0, \theta_1, \dots, \theta_{n_x}}_b]$



3. Logistic Regression Cost function

loss function:

一般用 squared error 表示 loss function: $J(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$

但在 LR 中很少使用, 因为它是 non-convex, 容易陷入 local optima

LR 的 loss function 为 $J(\hat{y}, y) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$

· if $y=1$, $J(\hat{y}, y) = -\log \hat{y}$ ← \hat{y} 越大越好

· if $y=0$, $J(\hat{y}, y) = -\log (1-\hat{y})$ ← \hat{y} 越小越好

· 目标是最小化 loss function

cost function

全部样本的 loss function → cost function

$$J(w, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log (1-\hat{y}^{(i)})]$$

4. Gradient Descent

$$w := w - \alpha \frac{\frac{dJ(w)}{dw}}{\left(\frac{\partial J(w, b)}{\partial w}\right)}$$

$$b := b - \alpha \frac{\frac{dJ(w)}{db}}{\left(\frac{\partial J(w, b)}{\partial b}\right)}$$

5. Gradient Descent in Logistic Regression

对于一个样本: $Z = w^T x + b$

$$\hat{y} = a = g(Z)$$

$$L(a, y) = -(y \log a + (1-y) \log(1-a))$$

$$\begin{aligned} & x_1 \\ & \downarrow w_1 \\ & x_2 \quad \rightarrow Z = w_1 x_1 + w_2 x_2 + b \rightarrow a = g(Z) \rightarrow L(a, y) \\ & \downarrow w_2 \\ & b \end{aligned}$$

对 dZ 求导:

$$dZ = \frac{\partial L}{\partial Z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial Z} = \left(-\frac{y}{a} + \frac{1-y}{1-a}\right) \cdot a(1-a) = a - y$$

对 w_1, w_2 和 b 求导:

$$dw_1 = \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial Z} \cdot \frac{\partial Z}{\partial w_1} = (a - y) x_1$$

$$dw_2 = (a - y) x_2, db = a - y$$

更新参数

$$w_1 := w_1 - \alpha dw_1, \quad w_2 := w_2 - \alpha dw_2, \quad b := b - \alpha db$$

6. m 个样本的 Gradient Descent

$$a^{(i)} = \hat{y}^{(i)} = g(Z^{(i)}) = g(w^T x^{(i)} + b)$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)}) = \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

$$dw_1 = \frac{1}{m} \sum_{i=1}^m x_1^{(i)} (a^{(i)} - y^{(i)}) \quad dw_2 = \frac{1}{m} \sum_{i=1}^m x_2^{(i)} (a^{(i)} - y^{(i)})$$

$$db = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

7. Vectorization

$$X: R^{n \times m}$$

$$W: R^{n \times 1}$$

$$b: R$$

$$Z: R^{1 \times m}$$

$$Z = W^T X + b$$

$$db = \frac{1}{m} \sum_{i=1}^m dZ^{(i)} = \frac{1}{m} np.sum(dZ)$$

$$dw = \frac{1}{m} X dZ^T$$

8. Logistic Regression cost function 的解釋

预测输出: $\hat{y} = g(W^T x + b)$, 其中 $g(z) = \frac{1}{1+e^{-z}}$

给定预测样本 x 下 $y=1$ 的概率 $\hat{y} = P(y=1|x)$

$$\text{if } y=1: P(y|x) = \hat{y}$$

$$\text{if } y=0: P(y|x) = 1 - \hat{y}$$

整合以上情况

$$P(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$$

由于 log 函数严格单调, 上述等式两边取 log

$$\log(P(y|x)) = \log(\hat{y}^y (1-\hat{y})^{1-y}) = y \log \hat{y} + (1-y) \log(1-\hat{y})$$

我们希望 $P(y|x)$ 越大越好, loss function 越小越好

$$L(\hat{y}, y) = -\log(P(y|x)) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

假设样本是 IID 的, 我们希望结果的联合概率越大越好

$$\max \prod_{i=1}^m P(y^{(i)}|x^{(i)})$$

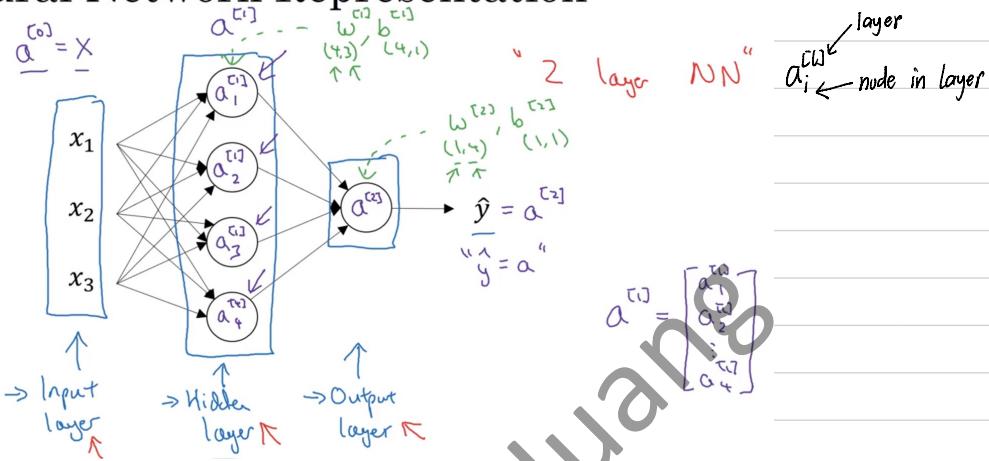
引入 log 函数, 上式简化

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

浅层神经网络

1. 简单双层 NN

Neural Network Representation



$$z_i^{(0)} = W^{(0)T} x + b_i^{(0)}, \quad a_i^{(0)} = g(z_i^{(0)})$$

$$z_2^{(0)} = W_2^{(0)T} x + b_2^{(0)}, \quad a_2^{(0)} = g(z_2^{(0)})$$

$$z_3^{(0)} = W_3^{(0)T} x + b_3^{(0)}, \quad a_3^{(0)} = g(z_3^{(0)})$$

$$z_4^{(0)} = W_4^{(0)T} x + b_4^{(0)}, \quad a_4^{(0)} = g(z_4^{(0)})$$

$$a^{(0)} = \begin{bmatrix} z_1^{(0)} \\ z_2^{(0)} \\ z_3^{(0)} \\ z_4^{(0)} \end{bmatrix} = g(z^{(0)})$$

$$\begin{bmatrix} W_1^{(1)} \\ W_2^{(1)} \\ W_3^{(1)} \\ W_4^{(1)} \end{bmatrix} \begin{bmatrix} x_1^{(0)} \\ x_2^{(0)} \\ x_3^{(0)} \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}$$

$$= \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}$$

为了得到最终结果，只需

$$\begin{aligned} z^{(1)} &= W^{(1)T} x + b^{(1)} \\ a^{(1)} &= g(z^{(1)}) \\ z^{(2)} &= W^{(2)T} a^{(1)} + b^{(2)} \\ a^{(2)} &= g(z^{(2)}) \end{aligned}$$

2. 多样本向量化

$$x \longrightarrow a^{(0)} = \hat{y}$$

$$x^{(0)} \longrightarrow a^{(0)(0)} = \hat{y}^{(0)} \longrightarrow$$

⋮

$$x^m \longrightarrow a^{(0)(m)} = \hat{y}^{(m)}$$

for $i = 1$ to m :

$$z^{(0)(i)} = W^{(0)T} x^{(i)} + b^{(0)}, \quad a^{(0)(i)} = g(z^{(0)(i)})$$

$$z^{(1)(i)} = W^{(1)T} a^{(0)(i)} + b^{(1)}, \quad a^{(1)(i)} = g(z^{(1)(i)})$$

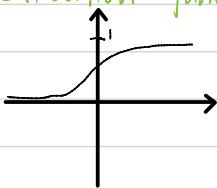
$$X = \begin{bmatrix} | & | & | & | \\ X^{(1)} & X^{(2)} & \dots & X^{(m)} \\ | & | & | & | \end{bmatrix} \quad Z^{(l)} = \begin{bmatrix} | & | & | \\ Z^{(l)(1)} & Z^{(l)(2)} & \dots & Z^{(l)(m)} \\ | & | & | \end{bmatrix} \quad W^{(l)} = \begin{bmatrix} | \\ | \\ | \end{bmatrix}$$

$$W^{(l)} X^{(1)} = \begin{bmatrix} | \\ Z^{(l)(1)} \\ | \end{bmatrix} \quad W^{(l)} X^{(2)} = \begin{bmatrix} | \\ Z^{(l)(2)} \\ | \end{bmatrix} \quad \dots \quad W^{(l)} X^{(m)} = \begin{bmatrix} | \\ Z^{(l)(m)} \\ | \end{bmatrix}$$

↓ ↓ ↓

$$W^{(l)} X = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \dots \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} | & | & | \\ Z^{(l)(1)} & Z^{(l)(2)} & \dots & Z^{(l)(m)} \\ | & | & | \end{bmatrix} = Z^{(l)}$$

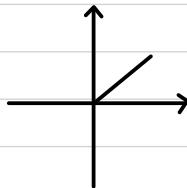
3. Activation functions



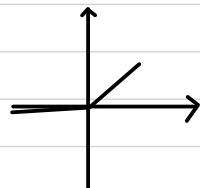
$$\text{sigmoid: } a = \frac{1}{1+e^{-z}} \\ a' = a(1-a)$$



$$\text{tanh: } a = \frac{e^z - e^{-z}}{e^z + e^{-z}} \\ a' = 1 - a^2$$



$$\text{ReLU: } a = \max(0, z)$$



$$\text{Leaky ReLU: } a = \max(0.01z, z)$$

Sigmoid: $[0, 1]$, 一般只用于二元分类

tanh: $[-1, 1]$, 表现优于 Sigmoid, 能起到归一化的作用

~~但是选择~~ ReLU: Sigmoid 与 tanh 在 $|z|$ 很大时, 梯度较小, 而 ReLU 改良了这点.

Leaky ReLU: $z < 0$ 时, 梯度下降仍不为 0

为什么要用非线性 Activation function?

如果不使用, 那么神经网络只是一个单纯的线性组合, 不如去掉隐藏层

4. NN 的梯度下降

参数: $W^{[0]}, b^{[0]}, W^{[1]}, b^{[1]}$

输入层特征数量: $N_x = n^{[0]}$

隐藏层单元个数: $n^{[1]}$

输出层单元个数: $n^{[2]} = 1$

$W^{[0]}$ 的维度 ($n^{[0]}, n^{[1]}$) $b^{[0]}$ 的维度 ($n^{[0]}, 1$)

$W^{[1]}$ 的维度 ($n^{[0]}, n^{[1]}$) $b^{[1]}$ 的维度 ($n^{[1]}, 1$)

Cost function: $J(W^{[0]}, b^{[0]}, W^{[1]}, b^{[1]}) = \frac{1}{m} \sum_{i=1}^m l(\hat{y}, y)$

Back Propagation: 逐元素乘积

$$dZ^{[0]} = A^{[0]} - Y \quad dW^{[0]} = \frac{1}{m} dZ^{[0]} A^{[0]T}$$

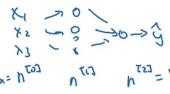
$$db^{[0]} = \frac{1}{m} \text{np.sum}(dZ^{[0]}, \text{axis}=1, \text{keepdims=True})$$

$$(n^{[0]}, m) \quad dZ^{[1]} = W^{[0]} dZ^{[0]} * g^{[0]'}(Z^{[0]}) \quad dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims=True})$$

Neural network gradients

$$\begin{aligned} x &\rightarrow z^{[1]} = W^{[1]}x + b^{[1]} \rightarrow a^{[1]} = \sigma(z^{[1]}) \rightarrow z^{[2]} = W^{[2]}x + b^{[2]} \rightarrow a^{[2]} = \sigma(z^{[2]}) \rightarrow L(a^{[2]}, y) \\ \frac{\partial L}{\partial w^{[1]}} &\rightarrow \frac{\partial L}{\partial z^{[2]}} = \frac{\partial L}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \rightarrow \frac{\partial L}{\partial z^{[1]}} = \frac{\partial L}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \rightarrow \frac{\partial L}{\partial w^{[2]}} = \frac{\partial L}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial w^{[2]}} \rightarrow " \frac{\partial L}{\partial w^{[1]}} = \frac{\partial L}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial w^{[1]}} " \\ \frac{\partial L}{\partial b^{[1]}} &\rightarrow \frac{\partial L}{\partial z^{[2]}} = \frac{\partial L}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \rightarrow \frac{\partial L}{\partial z^{[1]}} = \frac{\partial L}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \rightarrow \frac{\partial L}{\partial b^{[2]}} = \frac{\partial L}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial b^{[2]}} \rightarrow " \frac{\partial L}{\partial b^{[1]}} = \frac{\partial L}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial b^{[1]}} " \\ \frac{\partial L}{\partial x} &\rightarrow \frac{\partial L}{\partial z^{[2]}} = \frac{\partial L}{\partial a^{[2]}} \frac{\partial a^{[2]}}{\partial z^{[2]}} \rightarrow \frac{\partial L}{\partial z^{[1]}} = \frac{\partial L}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \rightarrow \frac{\partial L}{\partial x} = \frac{\partial L}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial x} \end{aligned}$$



Andrew Ng

5. 随机初始化

如果将神经元的 W 全部设置相同，那么每个神经元所计算出的梯度也相同。因

此无论经过多少次迭代，它们仍然是对称的。那么这将毫无意义。

$$W^{[0]} = \text{np. random. randn}((2, 2)) \times 0.01$$

$$b^{[0]} = \text{np. zeros}((2, 1))$$

乘 0.01 是因为如果使用 Sigmoid / tanh，使用较小的 W 值可以使 δ 接近于 0 附近，得到较大的梯度，加快收敛速度。

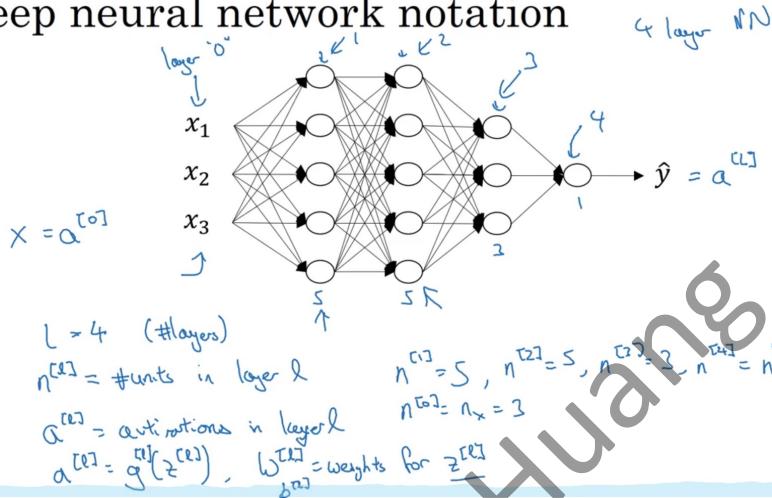
使用 ReLU / Leaky ReLU 则不存在这一问题。

Cruise Huang

深层神经网络

1. DNN的符号

Deep neural network notation



Andrew Ng

2. 计算矩阵的维度

$$\left. \begin{array}{l} W^{[l]} : (n^{[l]}, n^{[l-1]}) \\ dW^{[l]} : (n^{[l]}, n^{[l-1]}) \\ Z^{[l]} : (n^{[l]}, 1) \\ A^{[l]} : (n^{[l]}, 1) \end{array} \right\} \text{对于单个样本}$$

$\begin{matrix} & b^{[l]} : (n^{[l]}, 1) \\ & db^{[l]} : (n^{[l]}, 1) \end{matrix}$

对于单个样本: $A^{[l]} = g(Z^{[l]}) = W^{[l]} A^{[l-1]} + b^{[l]}$

对于多个样本: $A^{[l]} = g(Z^{[l]}) = W^{[l]} A^{[l-1]} + b^{[l]}$ 使用广播

3. 为什么使用深层表示

对于人脸识别，不同层的神经网络可以提取或者组合不同的特征，可以实现从局部到整体，由简单到复杂的模型学习。

对于语音识别，不同层可以分别学习音调、音素、单词等信息。

对于电路逻辑运算

Circuit theory and deep learning

Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

$y = x_1 \text{XOR } x_2 \text{XOR } x_3 \text{XOR } \dots \text{XOR } x_n$

$O(\log n)$

$x_1, x_2, x_3, \dots, x_n$

y

$O(2^n)$

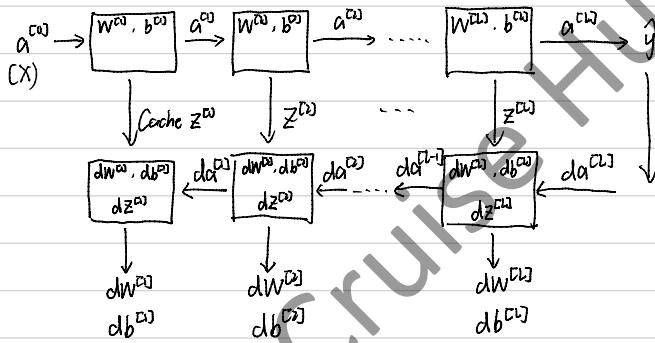
$N^{2^{n-1}}$ apparently large

Andrew Ng

$$\text{左边: } 1 + 2 + \dots + 2^{\log_2(n)-1} = n-1$$

$$\text{右边: } 2^{n-1}$$

4. 构建 DNN



5. 参数和超参数

参数: $W^{[0]}, b^{[0]}, W^{[1]}, \dots$

超参数: 控制参数的输出值

Learning rate: α

hidden Units: $n^{[0]}, n^{[1]}, \dots$

iteration: N

choice of activation function: Sigmoid, ReLU, tanh ...

hidden layer: L