

Cross function of Neural Network

定义: L : 网络的层数 $h_0(x)_k$ 代表输出的第 k 个类别的预测值
 S_l : 第 l 层的单元个数
 K : 输出层里单元的个数

回顾: logistic regression 中的 cost function (交叉熵).

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_0(x^{(i)})) + (1-y^{(i)}) \log(1-h_0(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

类似, 神经网络用在多分类时的 loss, 就是增加了一些权重求和.

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[y_k^{(i)} \log(h_0(x^{(i)})_k) + (1-y_k^{(i)}) \log(1-h_0(x^{(i)})_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\theta_{ji}^{(l)})^2$$

(PS: 用 sigmoid 函数将输出映射成概率, 值 $\in [0, 1]$, 但用于多分类的话, 其总和不一定为 1; softmax 是 0)
 类之间不排斥, 可能有多个正确答案.

第一部分: 对 output layer 的每个 cell 都计算 logistic loss 并求和 ($\sum \sum$)

第二部分: 对每个参数平方并求和 ($\sum \sum \sum$)
 before last layer
 layer

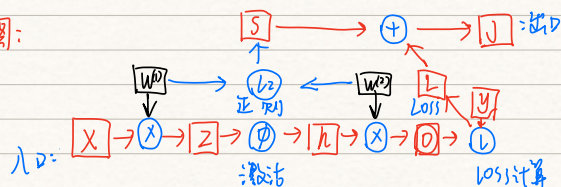
Backpropagation Algorithm

以前只知道“反向传播”这个词, 不知道它是一种算法

反向传播和 GD 用在线性回归中一样, 目标是最小化损失函数。也是一种参数优化 (寻找出一组最优的 θ_{ij}) 方法。
 逻辑 \rightarrow 指的是计算神经网络参数梯度的方法

Forward propagation: 正向传播 (FP): 按照 (从输入层到输出层) 计算和存储神经网络中每层的结果

计算图:



Backpropagation (反向传播, BP): 根据微积分中链式法则的规则, 从输出层到输入层的顺序

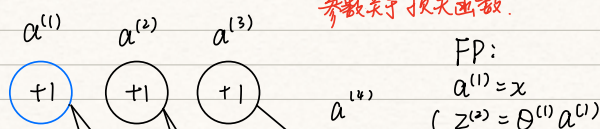
通过一次正向传播, 得到了预测值。

通过一次反向传播, 更新了参数 (一组 weight), 使得再一次 FP 得到的输出与真实值之间的误差减小。

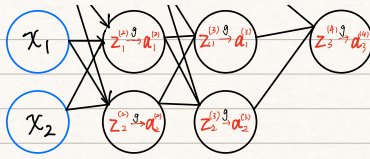
注意: “前向”: 从 input 层到 output 层
 “前一层”: 即: 是 l 的前一层

反向传播: 计算梯度, 梯度就是误差下降最快的方向。

参数关于损失函数。



BP: 从 output 层开始, 计算 $\delta^{(l)}$
 $\delta^{(l)} = a^{(l)} - y$ (即预测值与真实值之差)
 再利用 $\delta^{(l)} = ((\theta^{(l+1)})^T \cdot a^{(l)} \cdot x(1-a^{(l)}))$

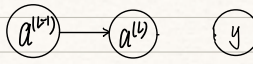


$$\begin{cases} a^{(2)} = g(z^{(2)}) \\ z^{(3)} = \theta^{(2)} a^{(2)} \\ a^{(3)} = g(z^{(3)}) \\ z^{(4)} = \theta^{(3)} a^{(3)} \\ a^{(4)} = g(z^{(4)}) = h_{\theta}(x) \end{cases}$$

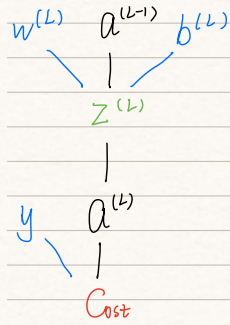
当层数多 后一层delta值 即 $g'(z^{(i)})$

与链式法则联系:

假设每层仅有一个结点:



每层有多个结点, 就是
(有好多条路, 加和符号!)



计算 cost 对 w 的微小变化有多敏感:

$$\frac{\partial \text{Cost}}{\partial w^{(L)}} = \frac{\partial \text{Cost}}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial z^{(L)}}{\partial w^{(L)}}$$

最后一层 上一层 上-层

