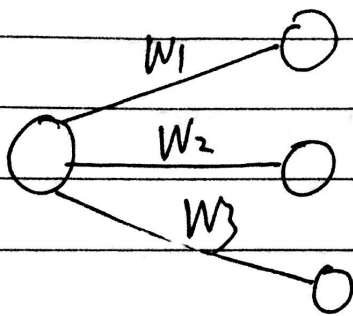


神经网络分类模型小结

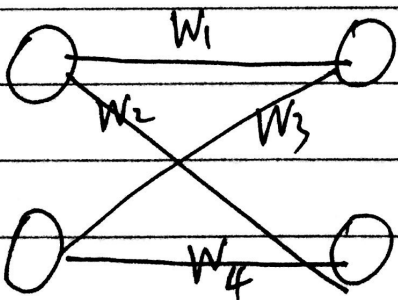
1. 首先需要设计好整体的网络架构。输入层为特征数，输出层为 label 的个数。重点是 hidden-layer 的层数以及每一层 unit 的个数。
2. 随机 initializer 各层参数矩阵 $\Theta^{(l)}$ 。为什么不能全部初始化为 0 呢？因为如果 $\Theta^{(l)}$ 全为 0，则每个 hidden-layer 的 activation 的值全部一样，最后 output units 全一样。虽然通过 BP 能更新这些 $\Theta^{(l)}$ ，但对于每个 neuron 的输出边



$$\frac{\partial J(0)}{\partial w_1} = \frac{\partial J(0)}{\partial w_2} = \frac{\partial J(0)}{\partial w_3}$$

这样便每次更新之后 $w_1 = w_2 = w_3$

当再进行 Forward Prop 时，由于每个 neuron



$$w_1 = w_2, w_3 = w_4$$

导致下一层每个 neuron 的输入值相同

进而输出值相同。

这样一个 neural network 模型的所有输出全是一致的。



无法取一个最大值将其划分为一个具体的 class

- 做 random initialization 的经验

$$\theta_{ij}^{(l)} \in [-\varepsilon, +\varepsilon) \quad \varepsilon = \frac{\sqrt{6}}{\sqrt{\text{input-unit} \times \text{output-unit}}}$$

input-unit: 第 l 层 neuron 个数

output-unit: 第 $(l+1)$ 层 neuron 个数

3. Forward Prop 计算当前 example 预测值。

第 l 层输出 $\times \theta^{(l)} \Rightarrow$ 第 $(l+1)$ 层输入

一层一层往后计算, 最后计算得出 output = $h(x)_k$

4. 计算 cost function. 将每个 example 的 $h(x)$ 计算最

后累加起来, 除上 m 。若忘正则化时, 加上所有 $\theta_{ij}^{(l)2}$ 注

意, 每个 bias 对应的参数不参与正则化。即忽略掉每个 $\theta^{(l)}$

的最后一列。

5. BP 来计算 $\frac{\partial J(x)}{\partial \theta_j}$ partial derivative

最后一层误差 $\delta^{(L)} = (y - \text{target})$ target 要转换为

binary vector. 然后向前计算上一层 $\delta^{(L-1)} = (\theta^{(L-1)})^T \cdot \delta^{(L)}$
 $\times \text{sigmoid}'(z^{(L-1)})$



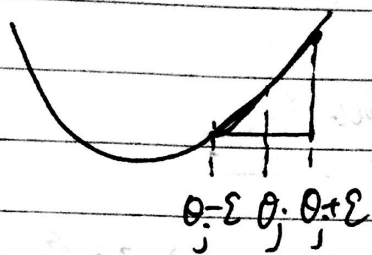
计算完每一层的 $\delta^{(l)}$ ，最后用公式

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} = \delta_i^{(l)} a_j^{(l)} \text{ 得到偏导数。将 } m \text{ 个 example 的}$$

偏导数和再除上 m 得到最终偏导数。注意，若总 J 则化时需要加上正则项的单独偏导，加上到总偏导中。

6. gradient checking。要保证通过公式计算的偏导与数值偏导基本相同。

将所 $\theta^{(l)}$ unroll 为 vector 将 $J(\theta)$ 看成



θ 的函数，然后其它 θ 不变， θ_j 改变 $\pm \epsilon$ 得

$$\text{需要近似相等于} \quad \leftarrow \quad \frac{f(\theta_j + \epsilon) - f(\theta_j - \epsilon)}{2\epsilon}$$

$$\frac{\partial J(\theta)}{\partial \theta_{ij}^{(l)}} \text{ 的值}$$

7. 已经使用 Forward 优化函数最小化 $J(\theta)$

已经计算出 J 代价 func 的值以及对每个参数的偏导，直接用 `fmincg` 来优化 cost function 使其值越来越小。最后得到最优 θ 的值

8. 用上述 θ 作预测。Forward Prop 然后取最大的值分为两类

