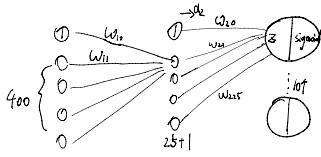$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \sum_{k=1}^{K} \left[ -y_k^{(i)} \log((h_\theta(x^{(i)}))_k) - (1-y_k^{(i)}) \log(1-(h_\theta(x^{(i)}))_k) \right] +$$
$$\frac{\lambda}{2m} \left[ \sum_{j=1}^{25} \sum_{k=1}^{400} (\Theta_{j,k}^{(1)})^2 + \sum_{j=1}^{10} \sum_{k=1}^{25} (\Theta_{j,k}^{(2)})^2 \right].$$

反向传播的过程  以 20x20 像素数字识别为例



$$J = -y \ln(h(x)) - (1-y)\ln(1-h(x))$$

以 $w_{21}$ 的偏导为例

$$\frac{\partial J}{\partial w_{21}} = \frac{-y}{h(x)} \frac{\partial h(x)}{\partial w_{21}} + \frac{(1-y)}{1-h(x)} \frac{\partial h(x)}{\partial w_{21}}$$

$$h(z) = sigmoid(z) = \frac{1}{1+e^{-z}} \qquad h'(z) = h(z)(1-h(z))$$

$$\frac{\partial J}{\partial w_{21}} = \frac{-y}{h(x)} h(x)(1-h(x)) \frac{\partial x}{\partial w_{21}} + \frac{(1-y)}{1-h(x)} h(x)(1-h(x)) \frac{\partial x}{\partial w_{21}}$$

$$= -y(1-h(x)) \frac{\partial x}{\partial w_{21}} + (1-y)h(x) \frac{\partial x}{\partial w_{21}}$$

$$= \left[ -y + yh(x) + h(x) - yh(x) \right] \frac{\partial x}{\partial w_{21}} = \left[ h(x) - y \right] \frac{\partial x}{\partial w_{21}} \longleftarrow$$

$\frac{\partial x}{\partial w_{21}}$ 其实就是 theta 中的 $w_{21}$ 对应的 $x$ ← 输入

这个就是隐藏层到输出层的梯度

---

然后基于上述结果继续求导，求出 输入层到隐藏层的梯度

例 $w_{11}$，由于 $w_{11}$ 输入隐藏层后，给输出层的各个输出了，

所以 求 $\frac{\partial J}{\partial w_{11}}$ 除了需要 $\frac{\partial J}{\partial w_{21}}$ 的，还需要其它的 $\frac{\partial J}{\partial w_{21}}$

这里先以 $\frac{\partial J}{\partial w_{21}}$ 求导为例

也是一个 sigmoid 的函数，同理求导

$$\frac{\partial J}{\partial w_{11}} = \left[ h(x) - y \right] \left( \frac{\partial a_{41}}{\partial w_{11}} \right) = \left[ h(x) - y \right] theta \cdot h(z)(1-h(z)) \cdot x$$