# 神经网络第五章

1. Train / Dev / Test sets
   训练, 验证, 测试
   Best output → 无偏估计
   (Train : Test = 7 : 3)
   (Train : Dev : Test = 6 : 2 : 2)
   Big data sample : (98% : 1% : 1%)
   　　　　　或 (99% : 0.5% : 0.5%)

2. Bias / Variance
   high bias → 欠拟合
   high variance → 过拟合 ...
   避免 high bias : 增加隐藏层个数,神经元个数
   　　　　　　训练时间, 换用更复杂的 NN
   避免 high variance : 增加训练样本数据, 或进行
   　　　　　　正则化 Regularization.

   (L2 regularization)
   $$J(w,b) = \frac{1}{m}\sum_{i=1}^{m}\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m}\|w\|_2^2$$
   $$\|w\|_2^2 = \sum_{j=1}^{n}w_j^2 = w^T w.$$

   (L1 regulation)
   $$J(w,b) = \frac{1}{m}\sum_{i=1}^{n}\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m}\|w\|_1$$
   $$\|w\|_1 = \sum_{j=1}^{n}|w_j|$$

   加入正则化项后
   $$dw^{[l]} = dW^{[l]}_{before} + \frac{\lambda}{m}W^{[l]}$$
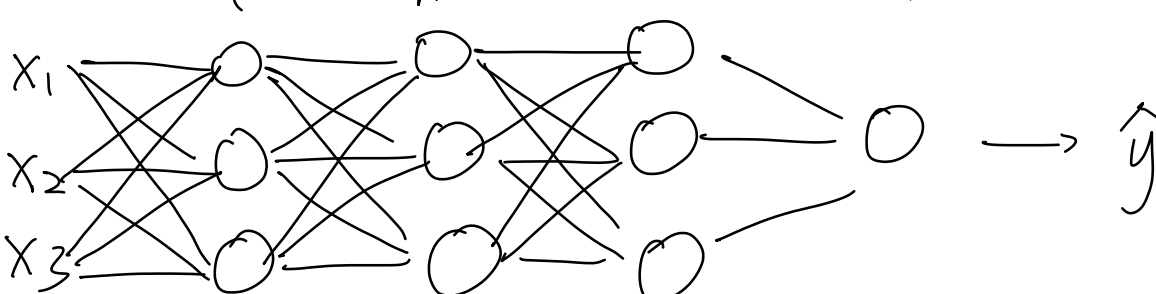   $$W^{[l]} := W^{[l]} - \alpha dw^{[l]}$$

   $L_2$ regularization → weight decay.
   $$W^{[l]} := W^{[l]} - \alpha dw^{[l]}$$
   $$= W^{[l]} - \alpha(dW^{[l]}_{before} + \frac{\lambda}{m}W^{[l]})$$
   $$= (1 - \alpha\frac{\lambda}{m})W^{[l]} - \alpha dw^{[l]}_{before}$$



dropout regularization :
　训练过程中, 对于每层神经元,按照一
　定概率, 排弃从神经网络中丢弃
　达到简化模型的效果,来避免过拟合

$dl = np.random.randn(al.shape[0], al.shape[1]) < keep\text{-}prob$
$al = np.multiply(al, dl)$
$al /= keep\text{-}prob.$
(设定 keep-prob = 0.8)

## Other regularization methods
① 对已有样本进行处理 → 更多样本
"猫"图片进行水平翻转,垂直翻转,
任意角度旋转,缩放,扩大 ......
(不需要增加额外成本, 却能防止过拟合...)
② 增加一些 noise
③ 防止过拟合, 提前 "early stopping"

## Normalizing inputs.
$$\mu = \frac{1}{m}\sum_{i=1}^{m}X^{(i)} \quad \sigma^2 = \frac{1}{m}\sum_{i=1}^{m}(X^{(i)})^2$$
$$X := \frac{x-\mu}{\sigma^2} \quad (对数据归一化)$$
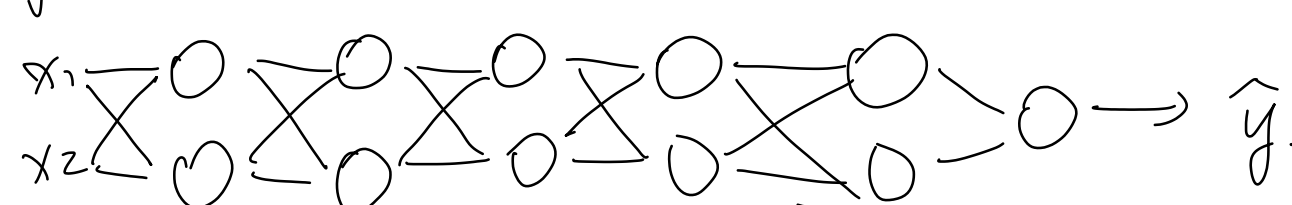$X_1 \in [1, 1000]$, $X_2 \in [0, 1]$
让输入归一化到同样的尺度上 ...
由于 $w_1, w_2$ 数值差异很大, 只能选择很小的学习因子
避免了振荡

## Vanishing and Exploding gradients
eg:



$$\hat{Y} = W^{[L]}W^{[L-1]}W^{[L-2]}\cdots W^{[3]}W^{[2]}W^{[1]}X$$

$\mathcal{L}$ 非常大时, 让步长过大或过小

## Weight Initialization for Deep Networks
· 改善梯度爆炸问题方法.
$$Z = w_1x_1 + w_2x_2 + \cdots + w_nx_n$$
$$a = g(z)$$
(让 w 与 n 有关, n越大, X越小)
$W^{[l]} = np.random.randn(n^{[l]}, n^{[l-1]}) * np.sqrt(1/n^{[l-1]})$
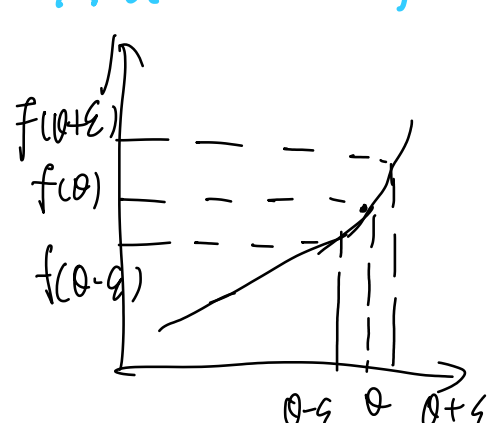(激活函数此时为 $tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$)
$w^{[l]} = np.random.randn(n^{[l]}, n^{[l-1]}) * np.sqrt(2/n^{[l-1]})$
(激... 为 ReLU)
$w^{[l]} = np.random.randn(n^{[l]}, n^{[l-1]}) * np.sqrt(2/n^{[l-1]} * n^{[l]})$

## Numerical approximation of gradients



$$g(\theta) = \frac{f(\theta+\varepsilon)-f(\theta-\varepsilon)}{2\varepsilon}$$
其中 $\varepsilon$ 足够小.

## Gradient checking
$$d\theta_{approx}[i] = \frac{J(\theta_1,\theta_2,\theta_i+\varepsilon,..)-J(\theta_1,\theta_2,..,\theta_i-\varepsilon,..)}{2\varepsilon}$$

$$\frac{\|d\theta_{approx}-d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2}$$
做不断逼近减小, 则反向梯度计算, 则无 bugs

## God checking Implementation
① 不要整个过程都用梯度, 仅作为 debug
② 不要忘掉正则项
③ 梯度检查时关闭 dropout, 完成后打开 dropout
④ 开始检查, 一段时间后再检查