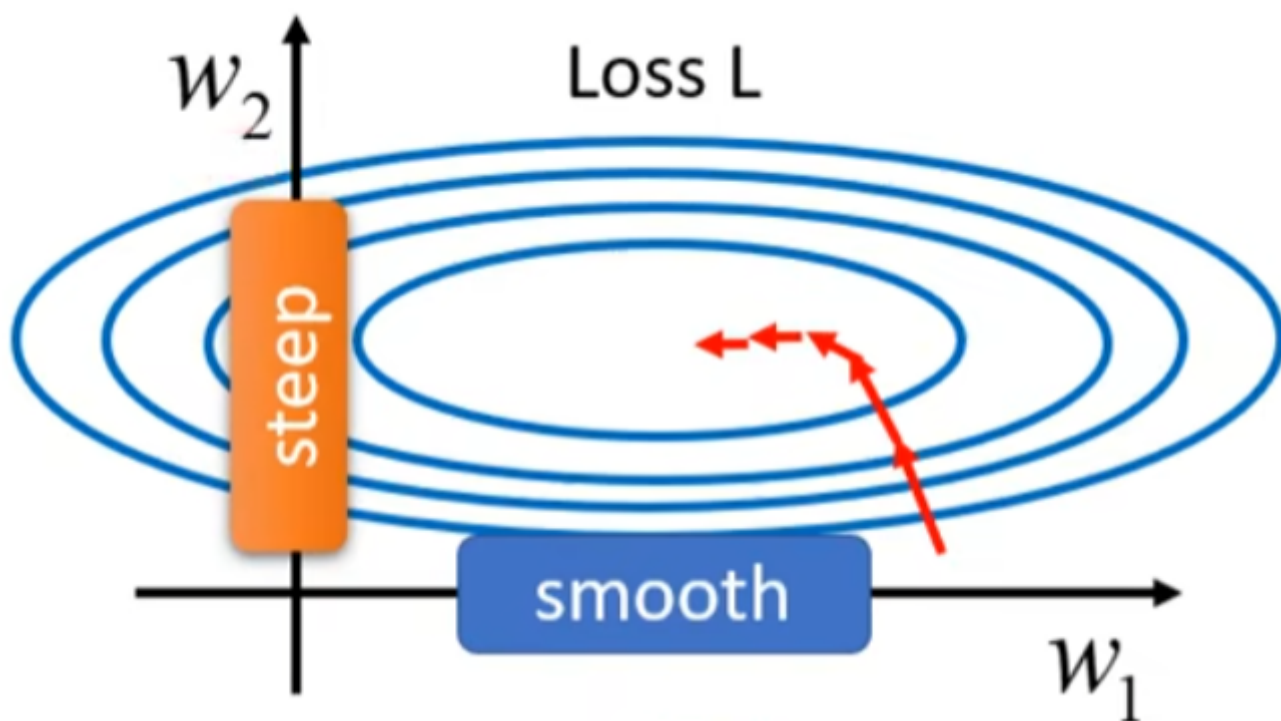


## 补充二：利用Batch Normalization优化训练

### 引入

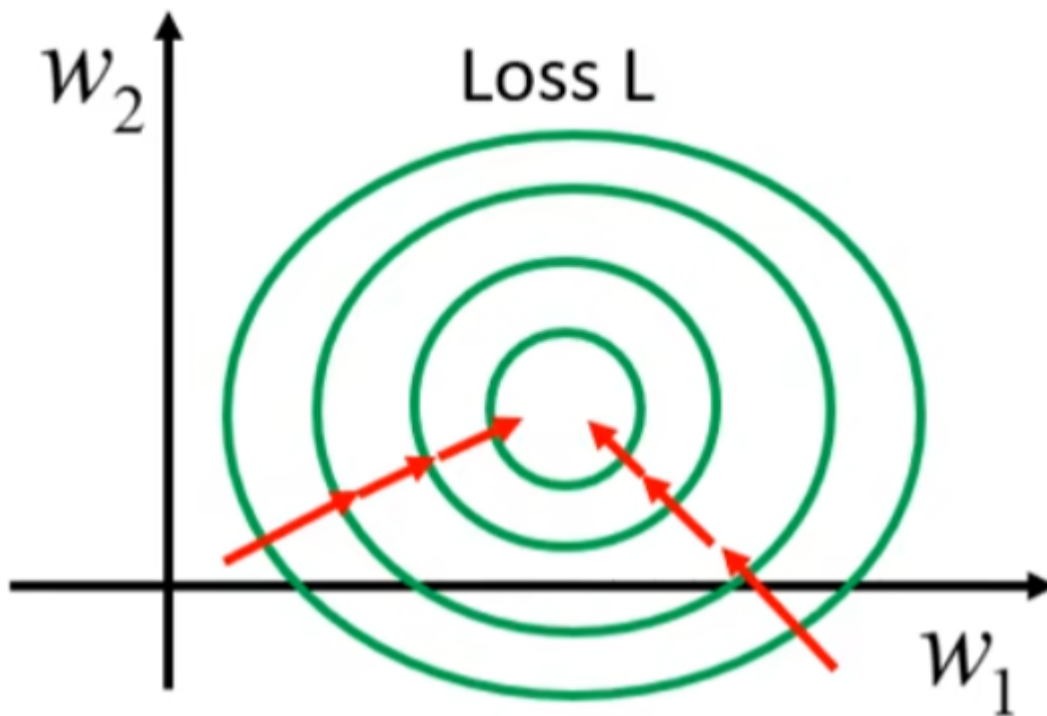
图中， $w_2$ 改变一点点，变化就很大，那能不能把error surface“铲平”（这样子就不会存在陡的时候需要慢慢摸索的情况了）



分析原因，由于 $y = wx + b$ ，真实值差距的 $e$ ：

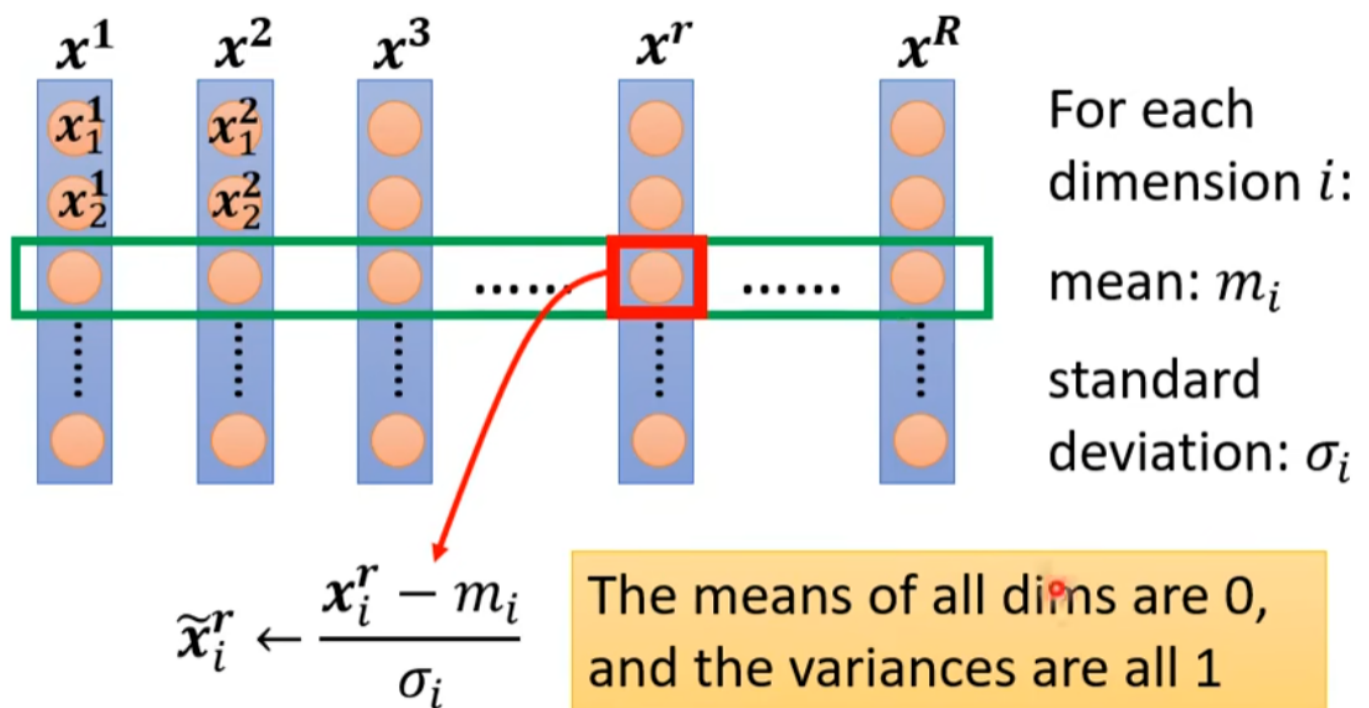
- 当 $x$ 很小的时候， $w$ 的变化对计算值 $y$ 的变化不大
- 当 $x$ 很大的时候，即使 $w$ 有很小的变化， $y$ 的值也会变化很大  $\rightarrow$   $e$ 的变化就很大

能不能让 $x$ 有相同的数据范围，让梯度下降尽量是一个“圆”：从任意一个点下降都很容易到谷底



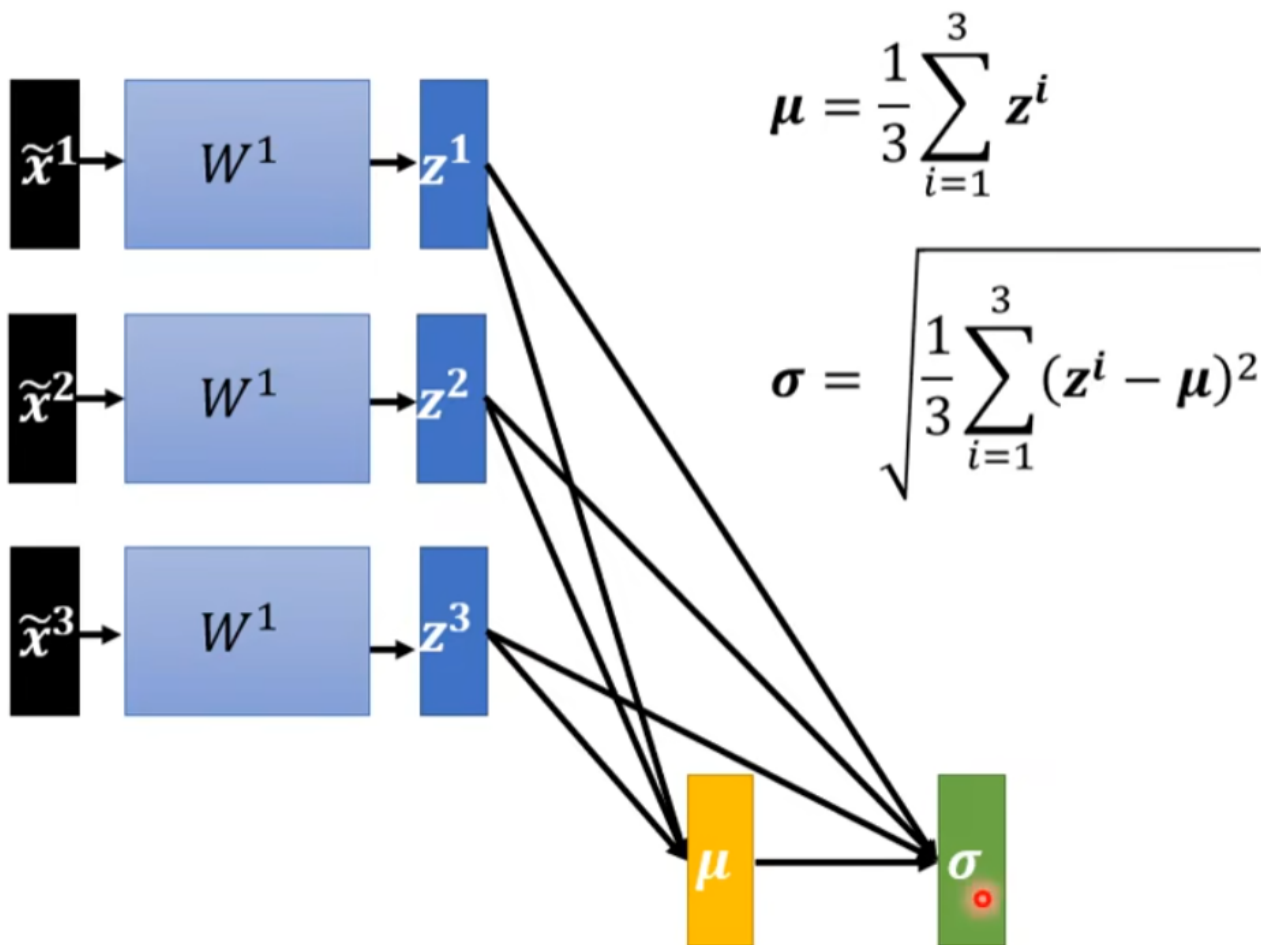
## Feature Normalization

(1) 直接对输入值 $x$ 归一化：让不同输入的相同位置做归一化，计算平均值和标准差。



此时，这一行的所有数值在0上下，偏差为1。

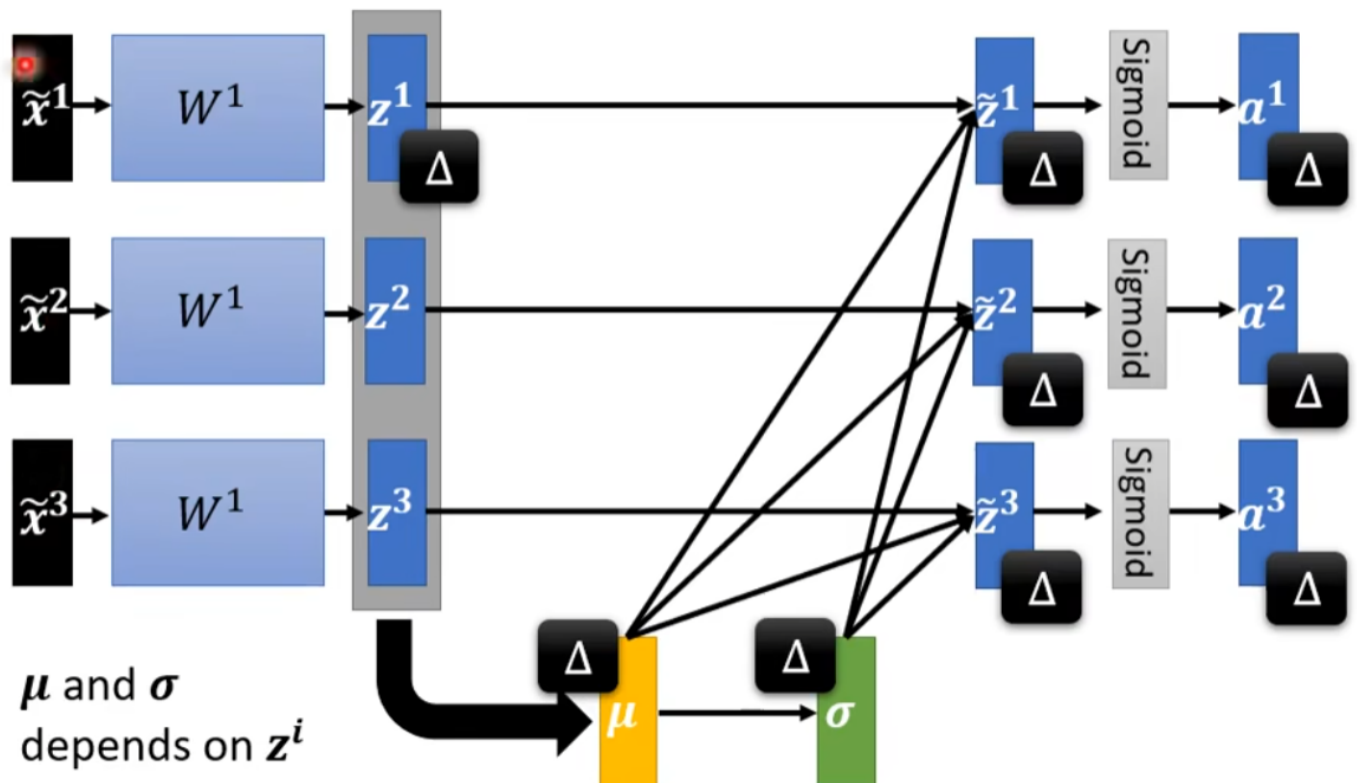
(2) 对中间的输出结果也需要归一化



但此时，一个数值的改变会影响 $\mu, \sigma$ ,从而到所有的输出变化：

## Considering Deep Learning

$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$



但在实际中，由于显存容量，无法算出整个数据集的归一化，因此只是对每个batch计算，即 batch normalization

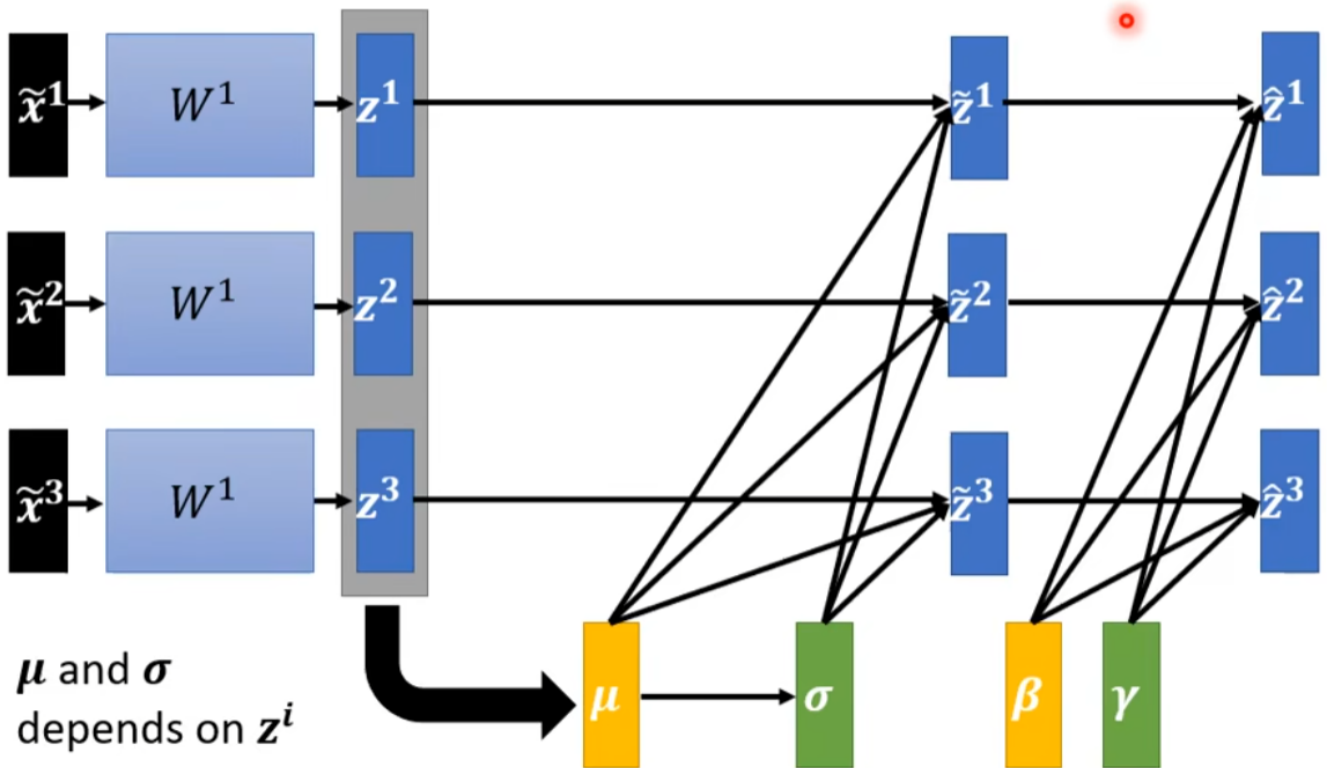
### Batch Normalization

但是在实际的训练过程中，需要引入 $\gamma, \beta$ ，让输出也存在一定的差异

# Batch normalization

$$\tilde{z}^i = \frac{z^i - \mu}{\sigma}$$

$$\hat{z}^i = \gamma \odot \tilde{z}^i + \beta$$



同时，在测试集中不一定存在batch，此时过程中之前固定batch数计算参数 $\mu, \sigma$ 的方法影响比较大。那么就可以利用平均数替代 $\bar{\mu} = p\bar{\mu} + (1 - p)p^t$

