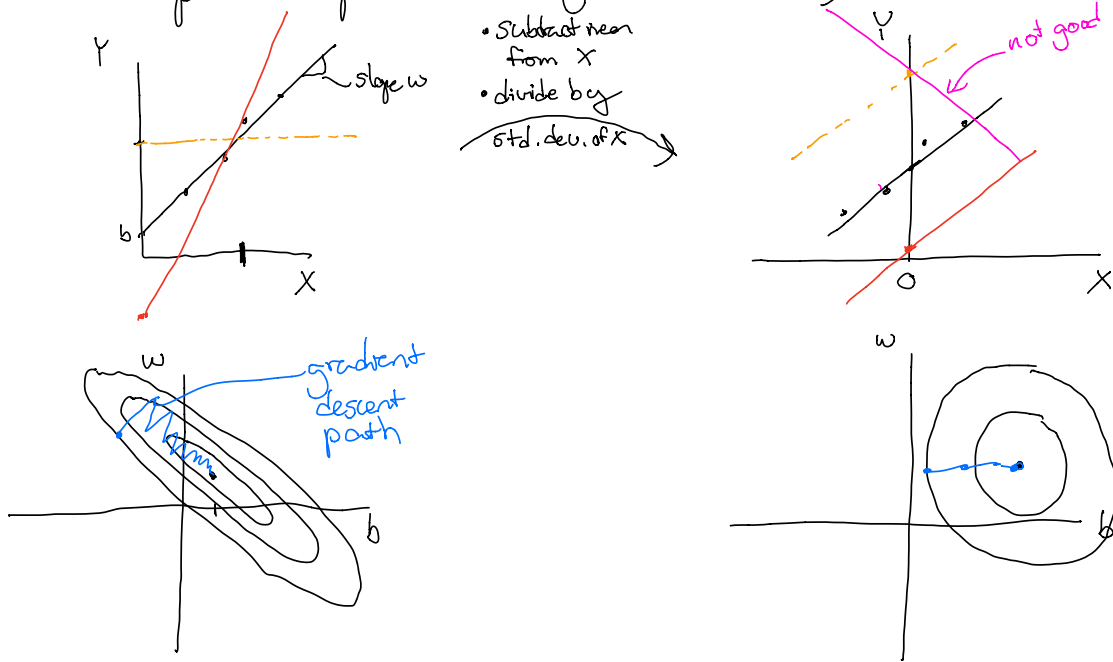
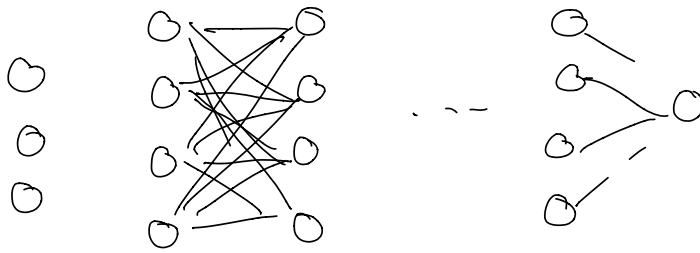


Reminder of normalizing network inputs:

Example: Simple linear regression: one  $X$ , numeric  $Y$



Larger network:



for this layer:

$$z^{[2]} = b^{[2]} + (w^{[2]})^T a^{[1]}$$

$$a^{[2]} = g^{[2]}(z^{[2]})$$

Idea of batch normalization:

- we standardize  $a^{[2]}$ , this fixes geometry of loss for  $b^{[2]}$  and  $w^{[2]}$ .

Normalizing inputs:

$$\mu^{[0]} = \frac{1}{m} \sum_{i=1}^m a^{(i)[0]} \quad \leftarrow \text{this is } X^{(i)}$$

$$\sigma^{2[0]} = \frac{1}{m-1} \sum_{i=1}^m (a^{(i)[0]} - \mu^{[0]})^2$$

$$\tilde{X}^{(i)} = \frac{X^{(i)} - \mu^{[0]}}{\sqrt{\sigma^{2[0]}}} \quad \left. \begin{array}{l} \text{average of } \tilde{X}^{(i)} \text{'s is } 0 \\ \text{std. dev. of } \tilde{X}^{(i)} \text{'s is } 1. \end{array} \right\}$$

Batch normalization for later layers:

$$\mu^{[l-1]} = \frac{1}{m} \sum_{i=1}^m a^{(i)[l-1]}$$

$$\sigma^{2[l-1]} = \frac{1}{m-1} \sum_{i=1}^m (a^{(i)[l-1]} - \mu^{[l-1]})^2$$

$$a_{\text{std}}^{(i)[l-1]} = \frac{a^{(i)[l-1]} - \mu^{[l-1]}}{\sqrt{\sigma^{2[l-1]} + \epsilon}} \quad \left. \begin{array}{l} \text{Average of } a_{\text{std}}^{(i)[l-1]} \text{'s is } 0 \\ \text{std dev. of } a_{\text{std}}^{(i)[l-1]} \text{ is } \approx 1 \end{array} \right\}$$

$\epsilon$  to prevent division by 0

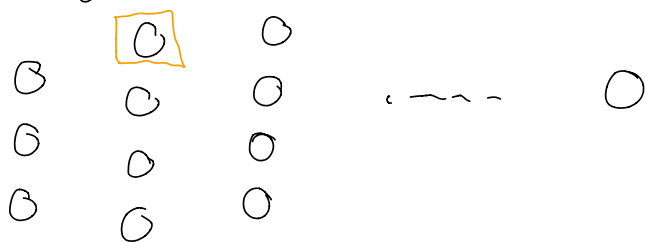
$$\tilde{a}^{(i)[l-1]} = \beta^{[l-1]} + \gamma^{[l-1]} * a_{\text{std}}^{(i)[l-1]} \quad \left. \begin{array}{l} \text{average of } \tilde{a}_j^{(i)[l-1]} \text{ is } \beta_j^{[l-1]} \\ \text{std. dev. of } \tilde{a}_j^{(i)[l-1]} \text{ is } \approx \sqrt{\gamma^{[l-1]}} \end{array} \right\}$$

shape  $(n_{l-1}, 1)$

these calculations occur within a Batch normalization layer  
the layer parameters are  $\beta$  and  $\gamma$ .

Keras function is `layers.BatchNormalization()`

A second motivation for Batch Normalization:  
helps reduce dependence among parameters in different layers of network



Suppose we have  $m=4$  obs.

$$a_1^{(1)[1]} = 2, \quad a_1^{(2)[1]} = 0, \quad a_1^{(3)[1]} = 2, \quad a_1^{(4)[1]} = 4$$

BN:  $\mu^{[1]}$  is a  $(4,1)$  shape array

$$\mu_1^{[1]} = \frac{2+0+2+4}{4} = 2$$

$$\sigma_1^{2[1]} = \frac{(2-2)^2 + (0-2)^2 + (2-2)^2 + (4-2)^2}{4} = \frac{8}{4} = 2$$

$$a_{1, std}^{(1)[1]} = \frac{2-2}{\sqrt{2}} \rightarrow 0 \quad (+\epsilon)$$

$$a_{1, std}^{(2)[1]} = \frac{0-2}{\sqrt{2}}$$

$$a_{1, std}^{(3)[1]} = \frac{2-2}{\sqrt{2}} \rightarrow 0$$

$$a_{1, std}^{(4)[1]} = \frac{4-2}{\sqrt{2}}$$

$$-\sqrt{2}$$

$$\sqrt{2}$$

$$\tilde{a}_1^{(1)[1]} = \beta \quad \tilde{a}_1^{(2)[1]} = \beta - \sqrt{2}\gamma \quad \tilde{a}_1^{(3)[1]} = \beta \quad \tilde{a}_1^{(4)[1]} = \beta + \sqrt{2}\gamma$$

What happens if  $w^{[1]}$  changes and  
 $a_1^{(1)[1]} = 200, \quad a_1^{(2)[1]} = 0, \quad a_1^{(3)[1]} = 200, \quad a_1^{(4)[1]} = 400.$

( after B.N.,

$$\tilde{a}_1^{(1)[1]} = \beta, \quad \tilde{a}_1^{(2)[1]} = \beta - \sqrt{2}\gamma, \quad \tilde{a}_1^{(3)[1]} = \beta, \quad \tilde{a}_1^{(4)[1]} = \beta + \sqrt{2}\gamma$$

we don't need a big adjustment to param's for layer 2 to account for a big change in activation outputs from layer 1.