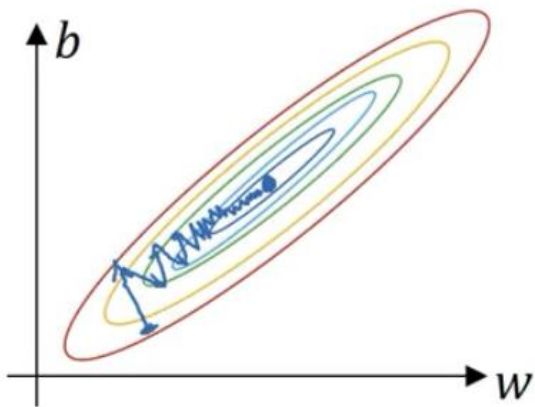
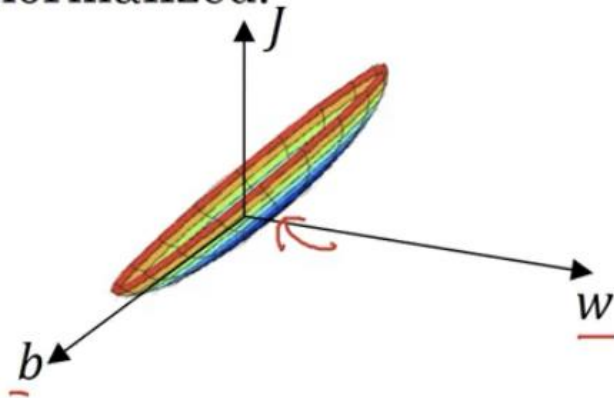


# Why normalize inputs?

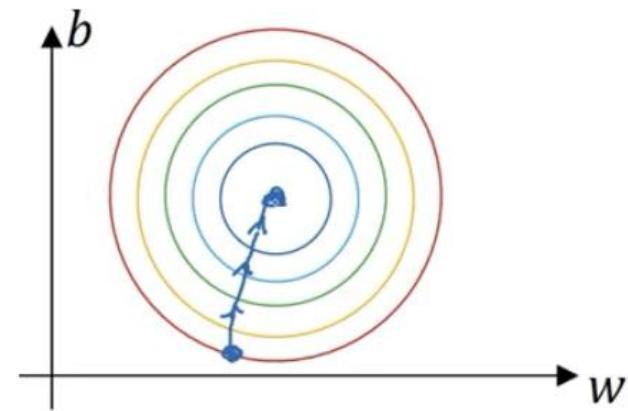
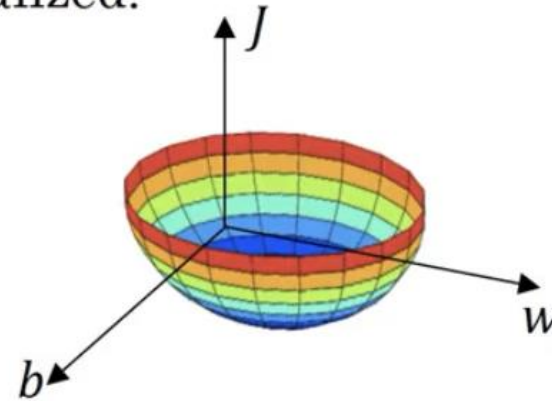
$w_1, x_1: 1 \dots 1000$   
 $w_2, x_2: 0 \dots 1$

Unnormalized:



$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Normalized:



# Hyperparameters To Tune

$\alpha$

$\beta$  0.9

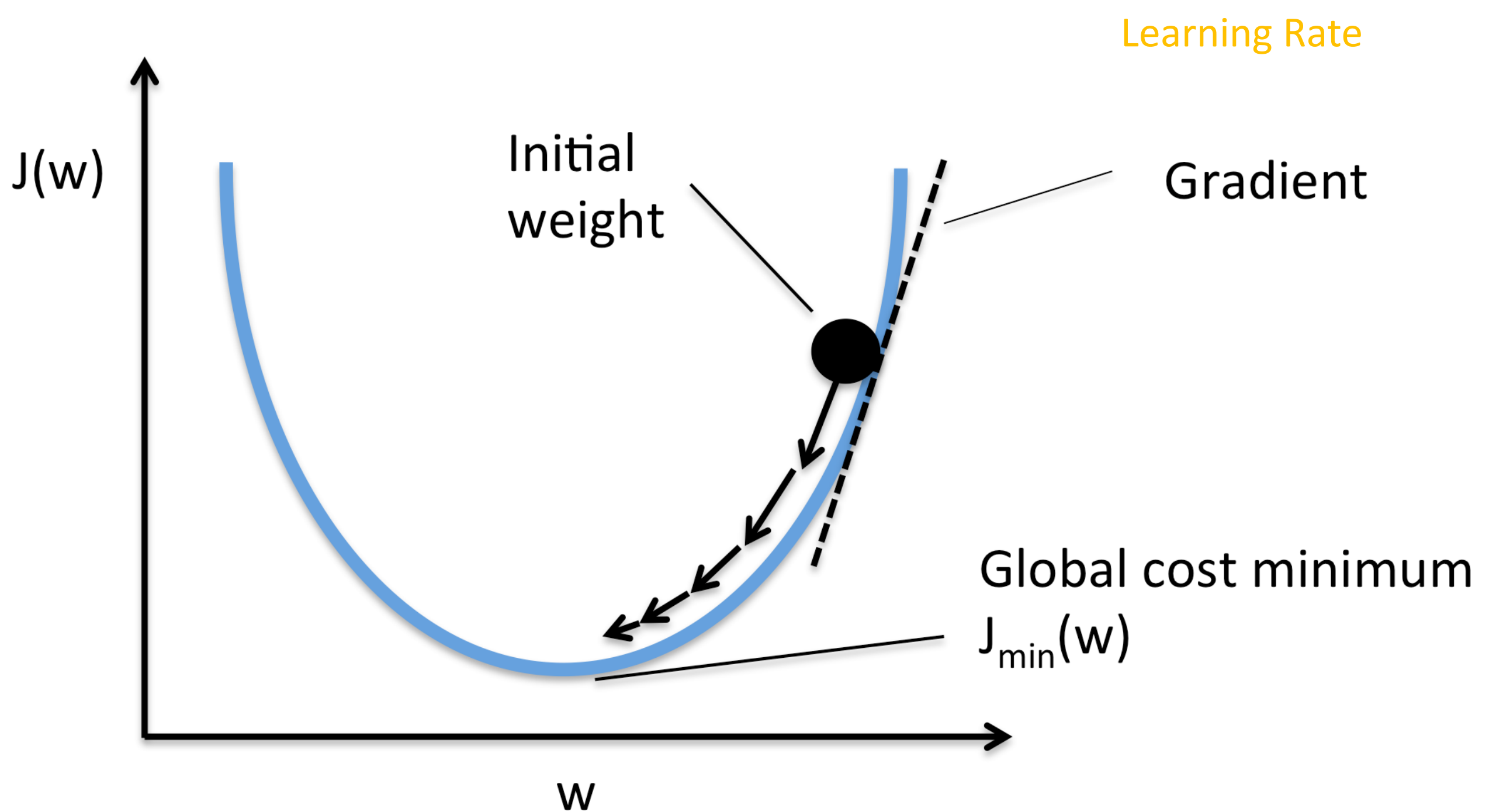
$\beta_1, \beta_2, \epsilon$

#layers

#hidden units

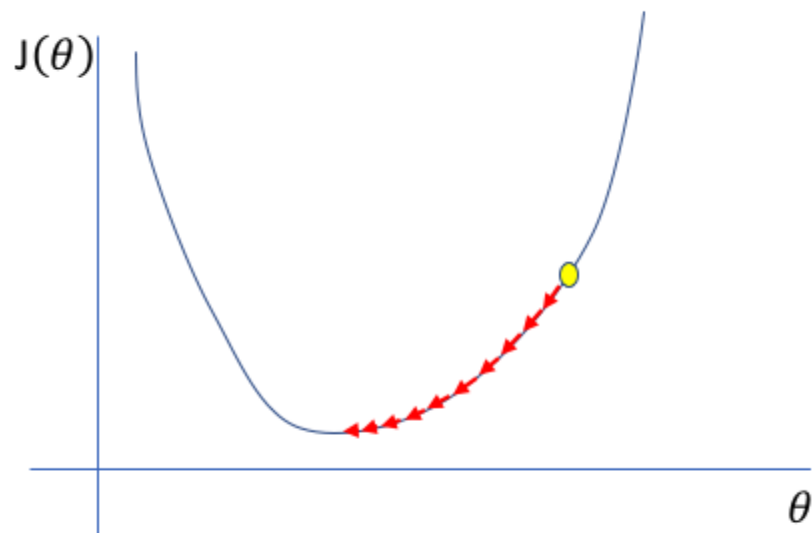
learning rate decay

mini-batch size



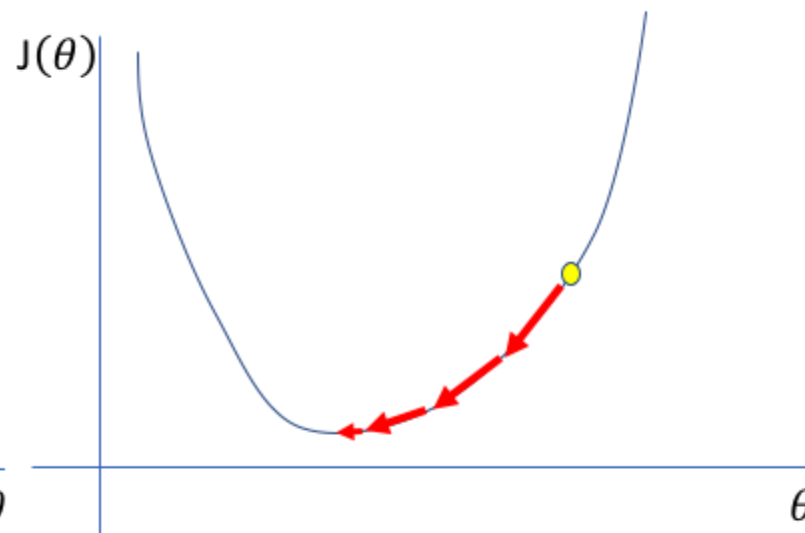
# Learning Rate

Too low



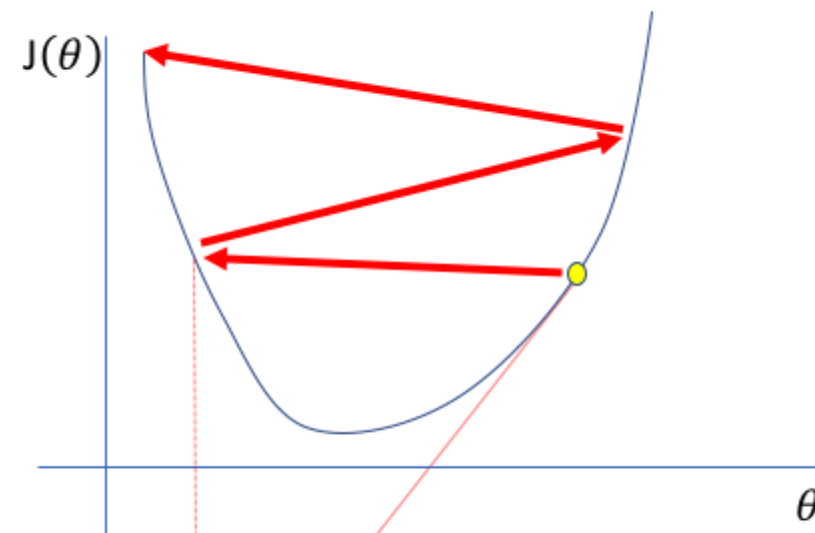
A small learning rate requires many updates before reaching the minimum point

Just right



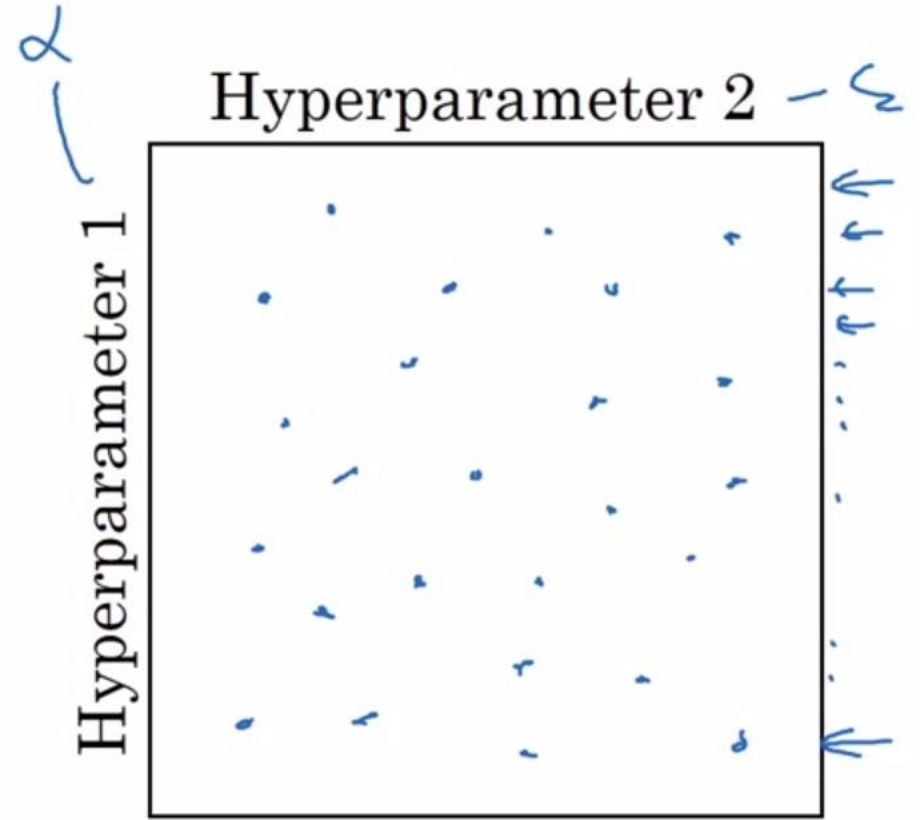
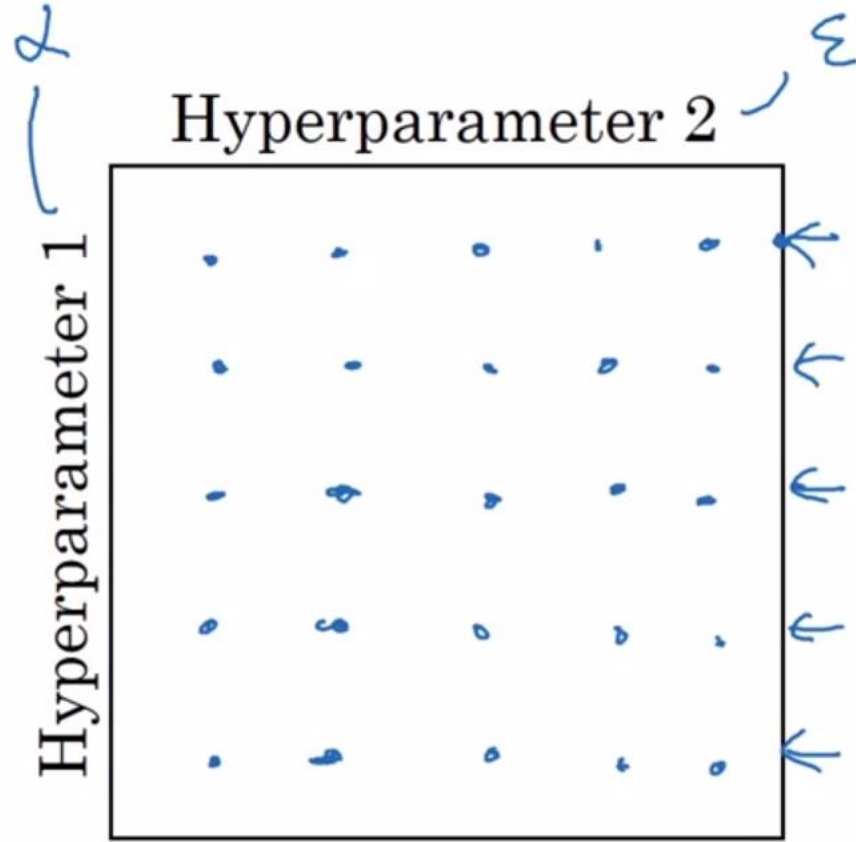
The optimal learning rate swiftly reaches the minimum point

Too high

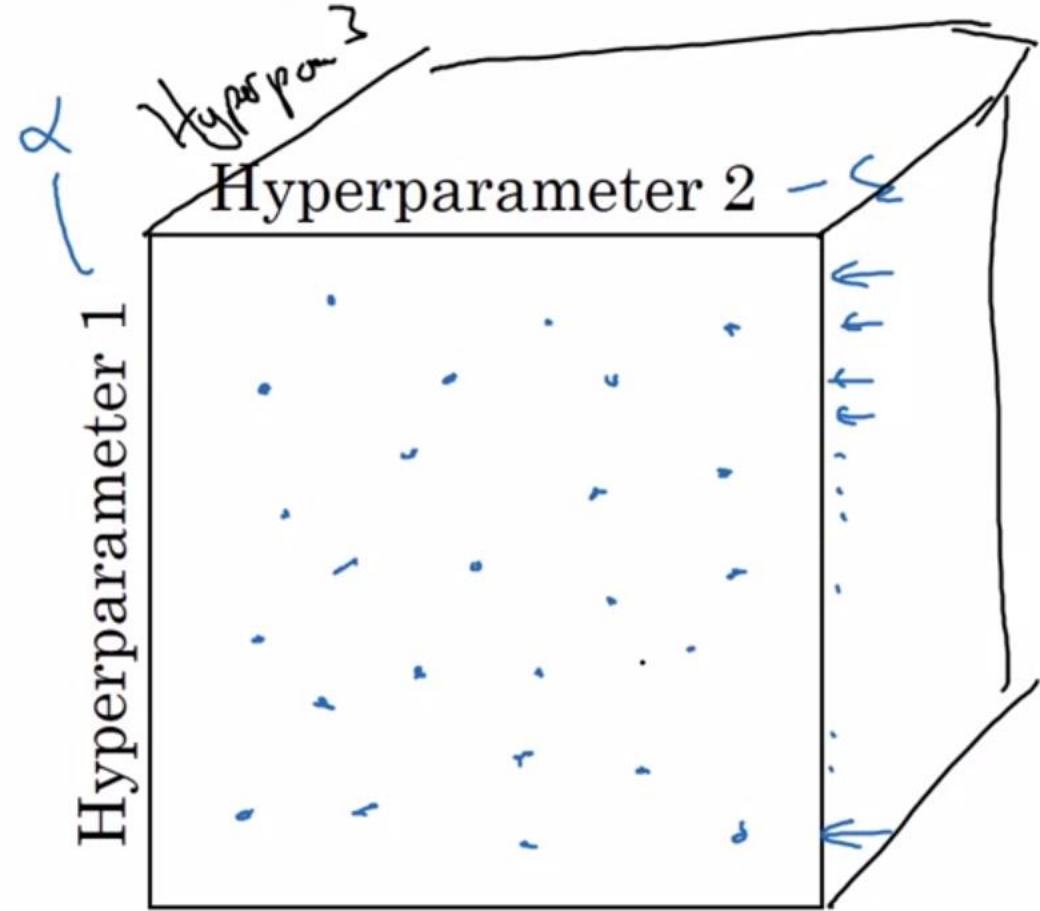
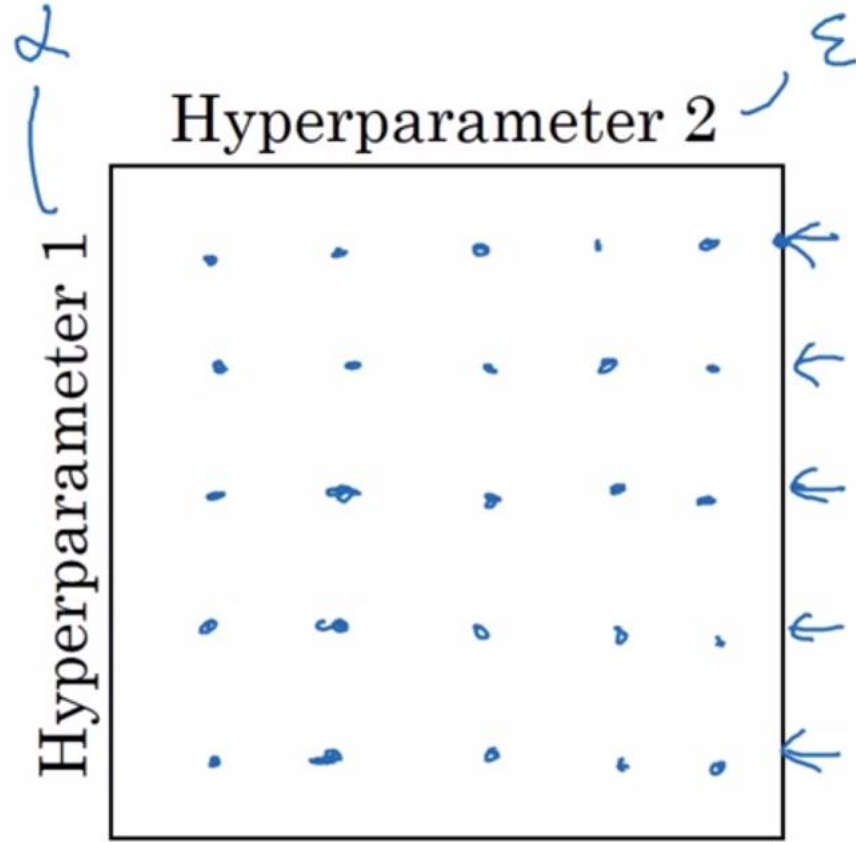


Too large of a learning rate causes drastic updates which lead to divergent behaviors

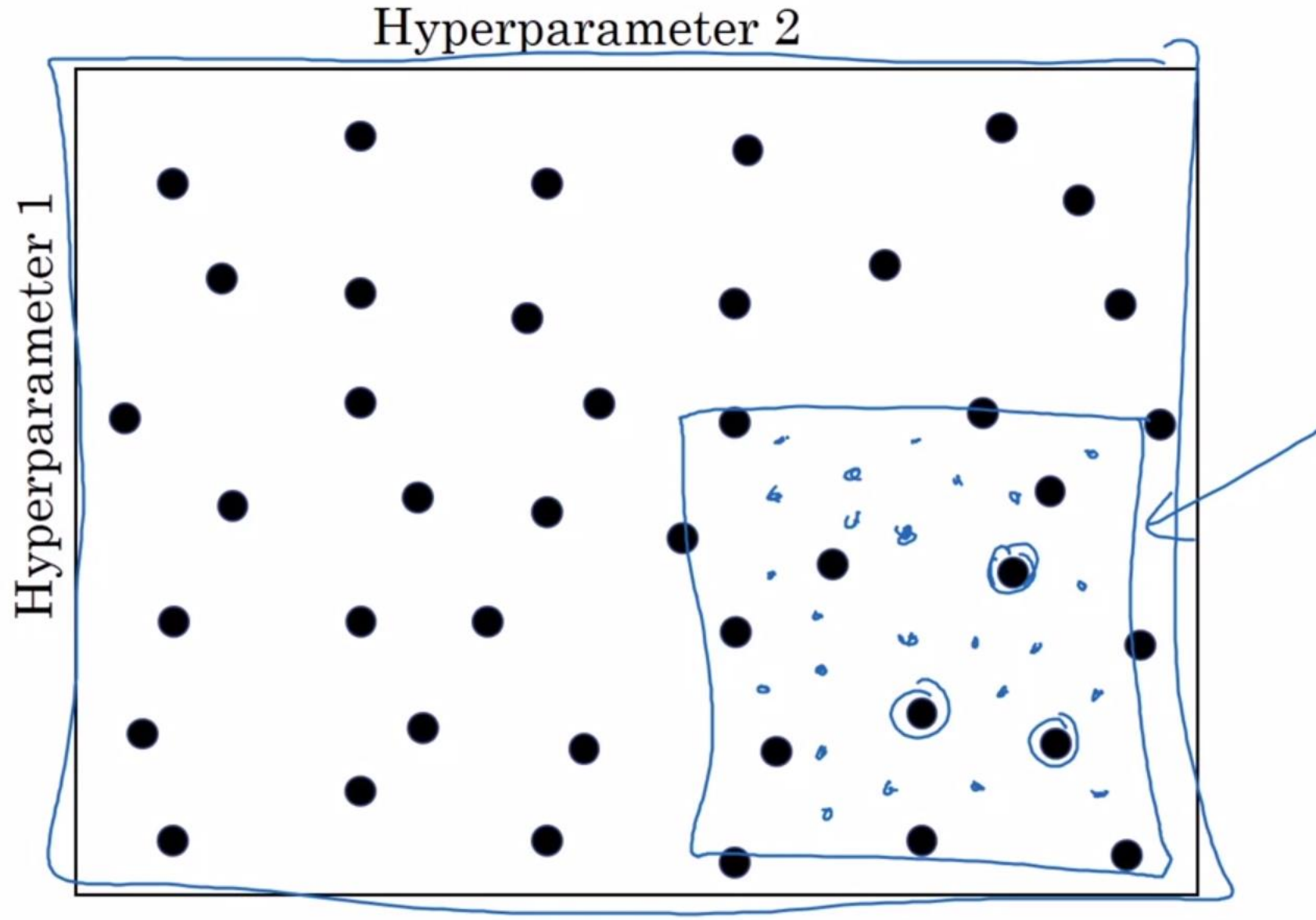
# Try random values: Don't use a grid



# Try random values: Don't use a grid



# Coarse to fine



# Picking hyperparameters at random

$$\rightarrow n^{\text{test}} = 50, \dots, 100$$



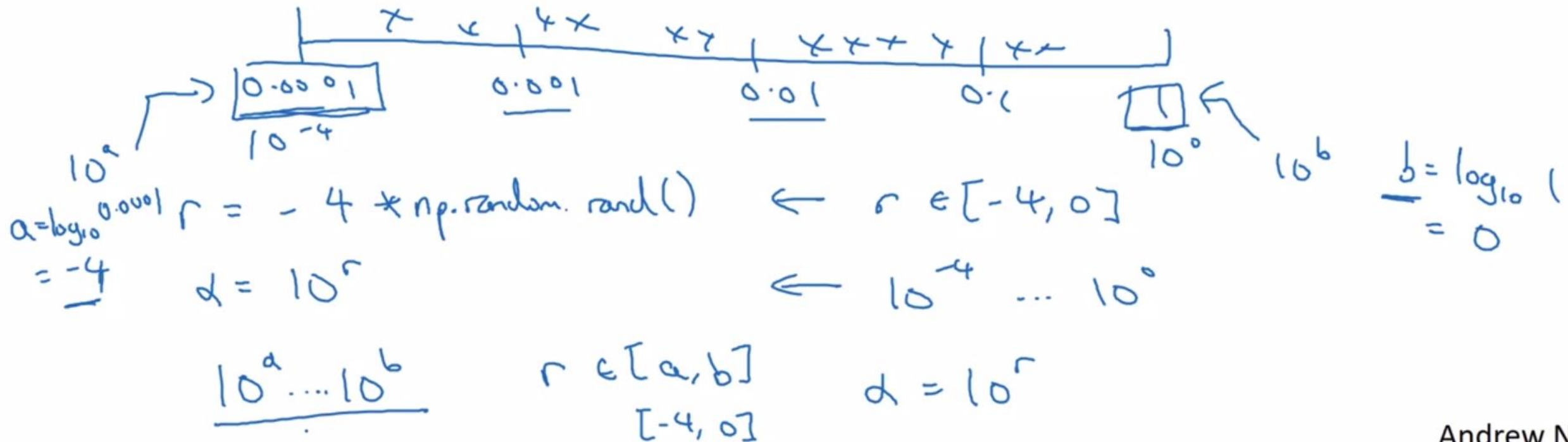
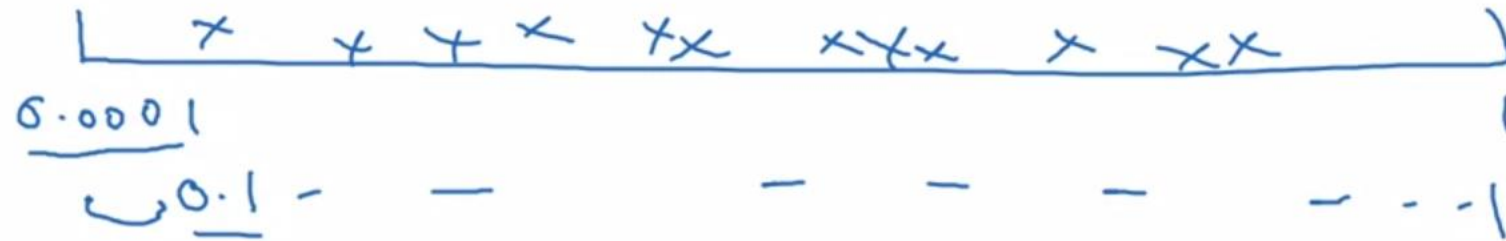
$$\rightarrow \# \text{layers} \quad L : 2 - 4$$

$$2, 3, 4$$

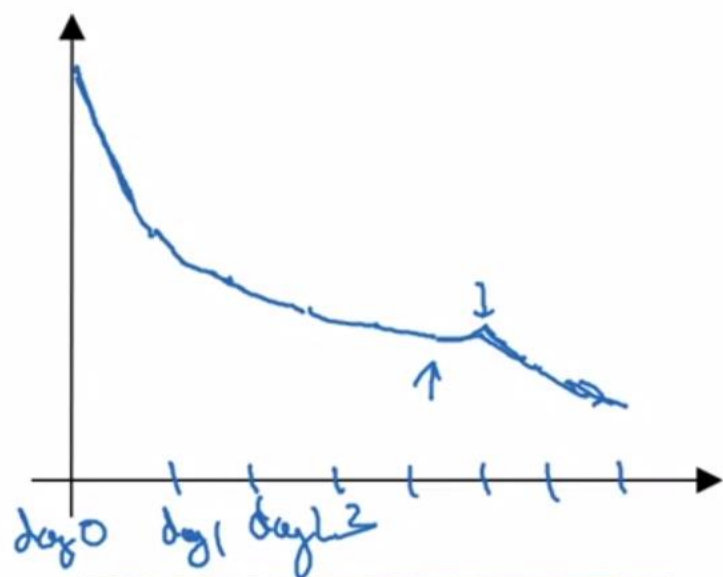


# Appropriate scale for hyperparameters

$$\alpha = 0.0001, \dots, 1$$

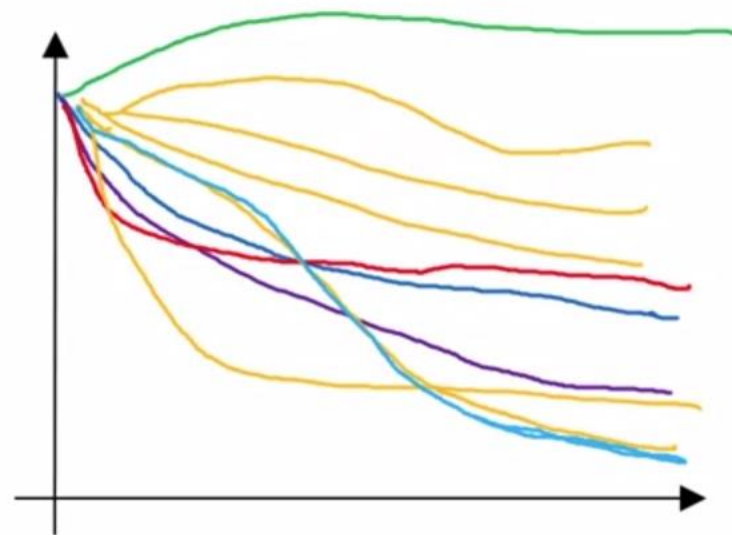


# Babysitting one model



Panda ←

# Training many models in parallel



Caviar ←

# Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On iteration  $t$ :

Compute  $dw, db$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \quad \leftarrow \text{"momentum"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) dw^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \quad \leftarrow \text{"RMSprop"} \beta_2$$

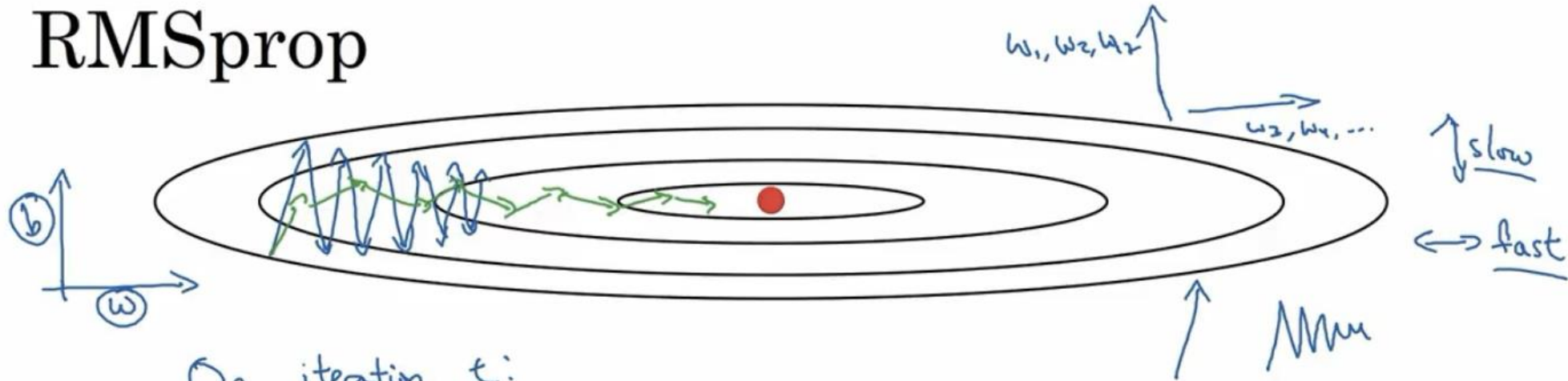
$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

$$W := W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}$$

$$b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

# RMSprop



On iteration  $t$ :

Compute  $dw, db$  on current mini-batch

$$\underline{S_{dw}} = \beta_2 S_{dw} + (1 - \beta_2) \underbrace{dw^2}_{\text{element-wise}} \leftarrow \text{small}$$

$$\rightarrow \underline{S_{db}} = \beta_2 S_{db} + (1 - \beta_2) \underline{db^2} \leftarrow \text{large}$$

$$w := w - \frac{\alpha dw}{\sqrt{S_{dw}}} \leftarrow$$

$$b := b - \frac{\alpha db}{\sqrt{S_{db}}} \leftarrow$$

# Implementation details

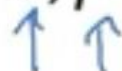
On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

Hyperparameters:  $\alpha, \beta$   


$$\underline{\beta = 0.9}$$