# Computational Intelligence

Samaneh Hosseini
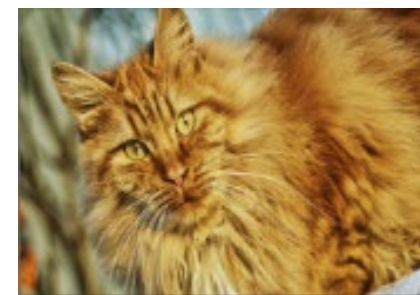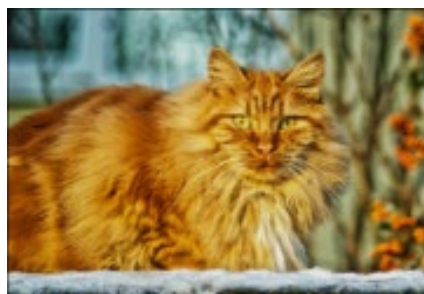
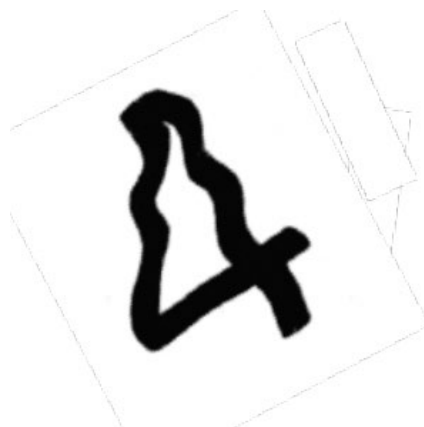Isfahan University of Technology

# Outline

- Other Regularization Techniques

  - Data Augmentation

  - Early Stopping

- Setting up Your Optimization Problem

  - Normalizing Inputs

  - Vanishing/Exploding Gradients

  - Numerical approximation of gradients

  - Gradient Checking

  - Gradient Checking Implementation Notes

# Other Regularization Techniques: Data Augmentation

# Data Augmentation

# Early Stopping

- Stop training before we have a chance to overfit



Loss

Training Iterations

dev set

train set

# Early Stopping

- Stop training before we have a chance to overfit



Loss / Training Iterations

**Legend**
dev set
train set

# Early Stopping

- Stop training before we have a chance to overfit



Legend

dev set

train set

# Early Stopping

- Stop training before we have a chance to overfit



Loss / Training Iterations

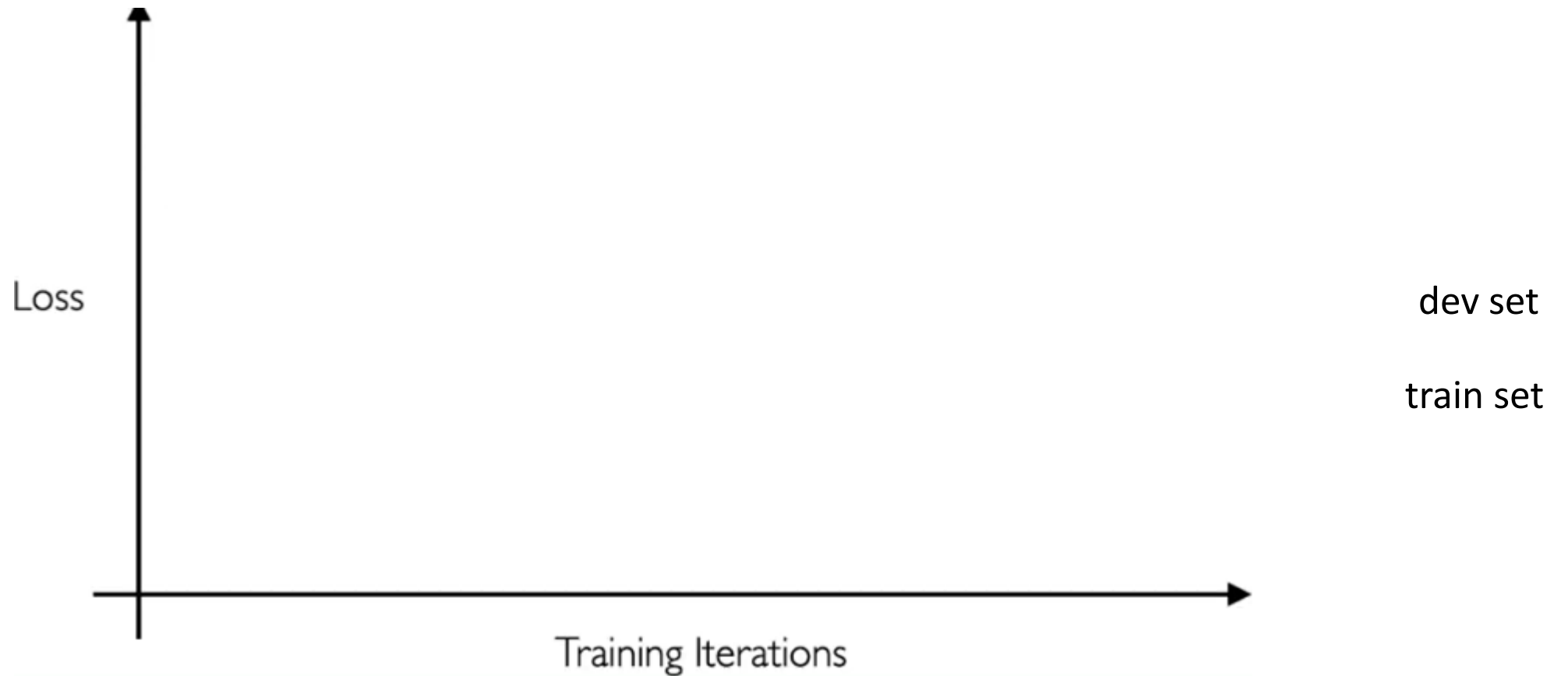**Legend**

dev set

train set

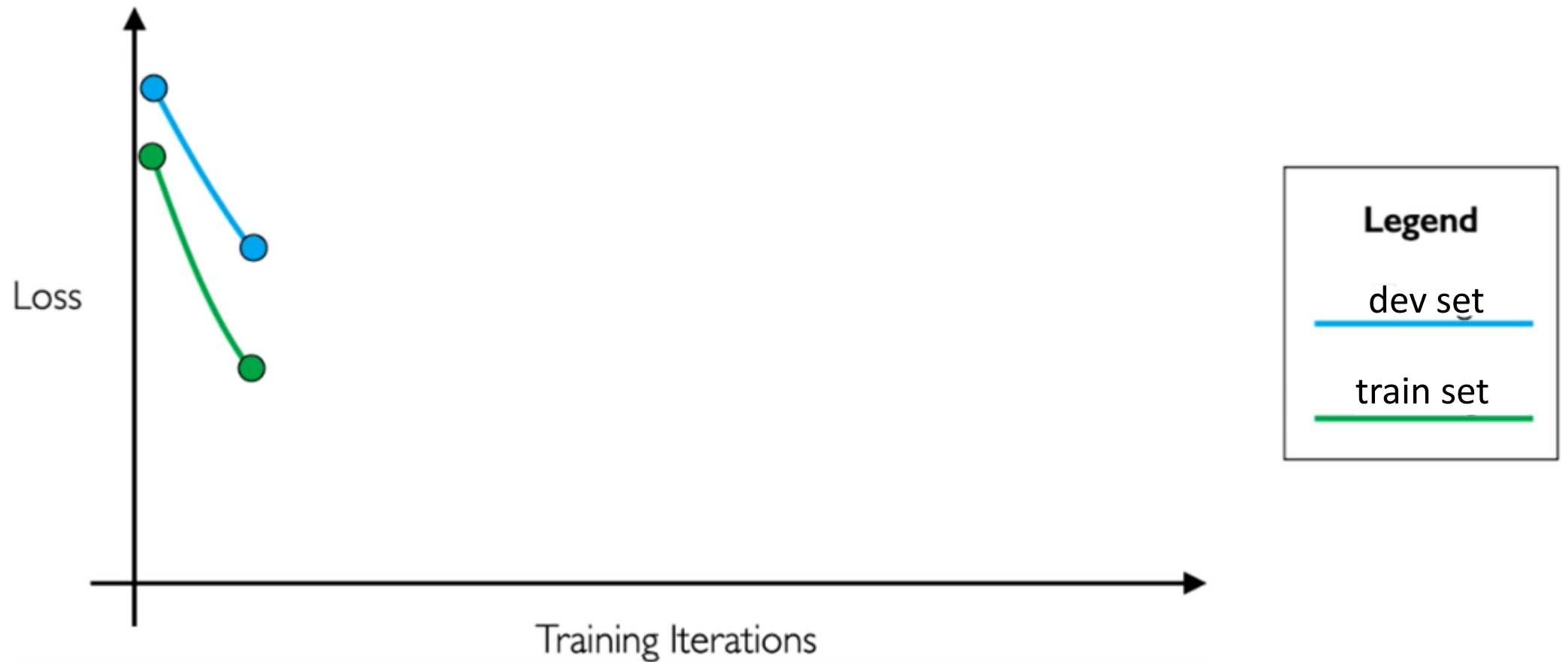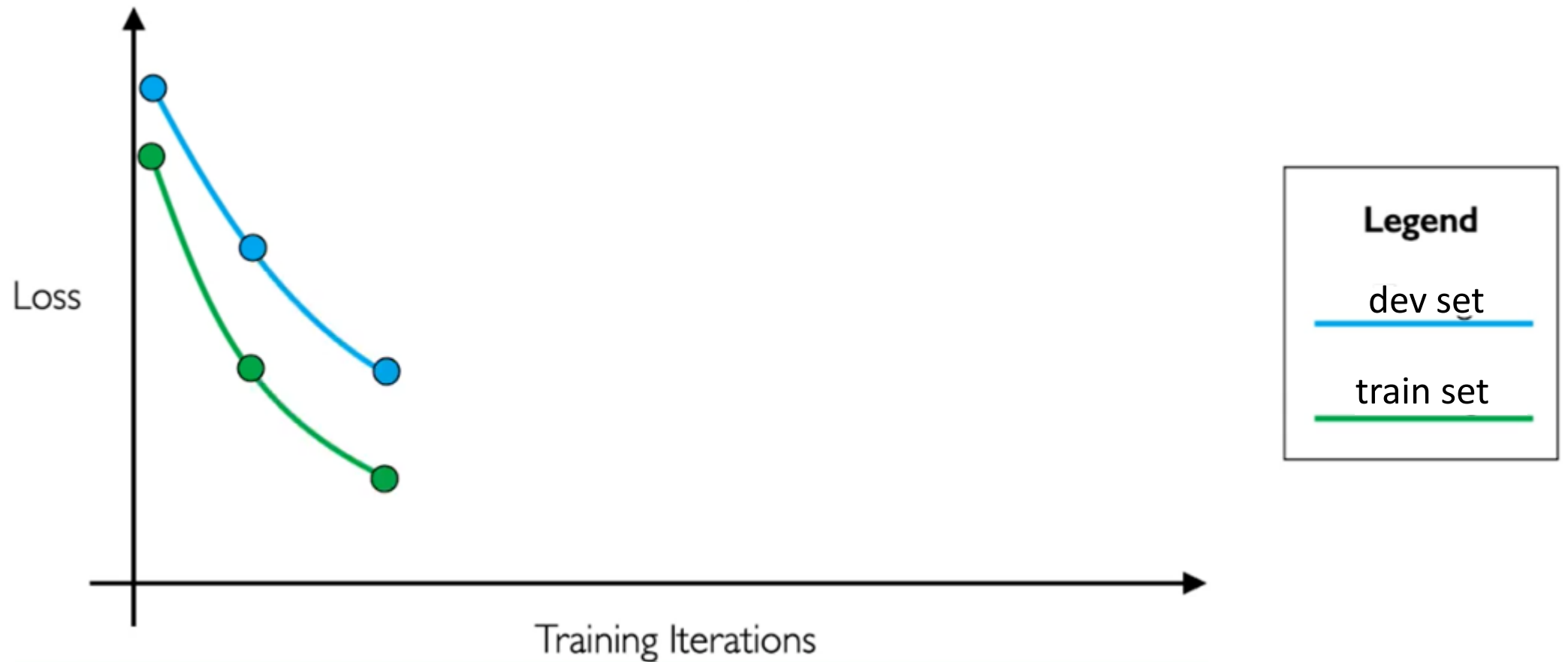# Early Stopping

- Stop training before we have a chance to overfit

# Early Stopping

- Stop training before we have a chance to overfit



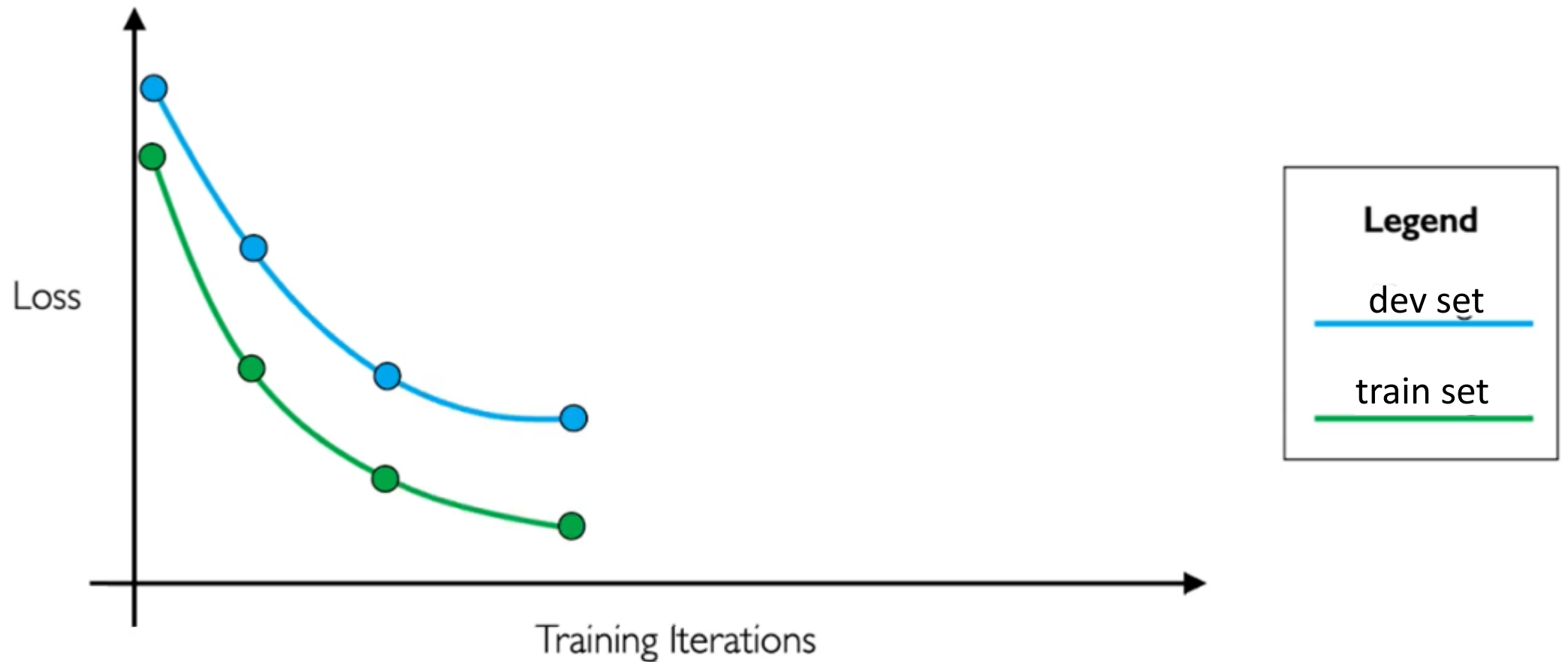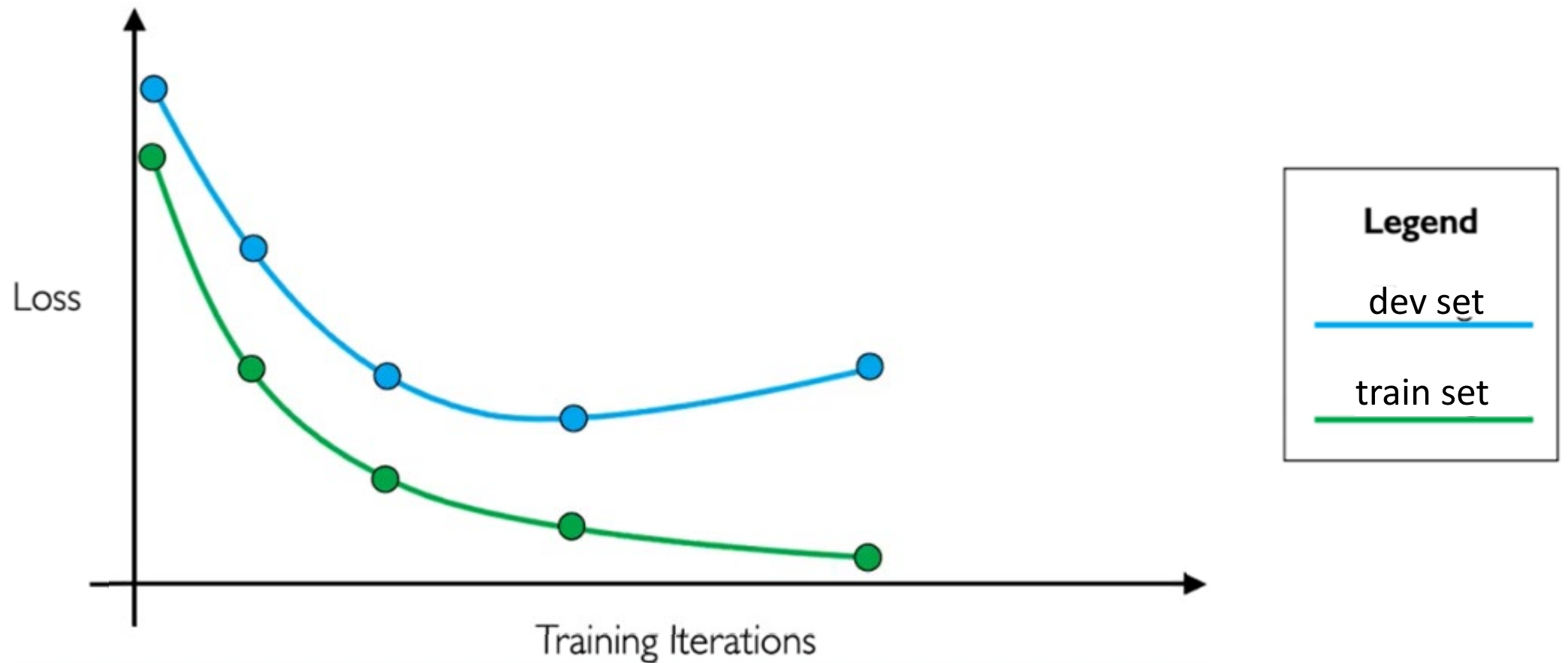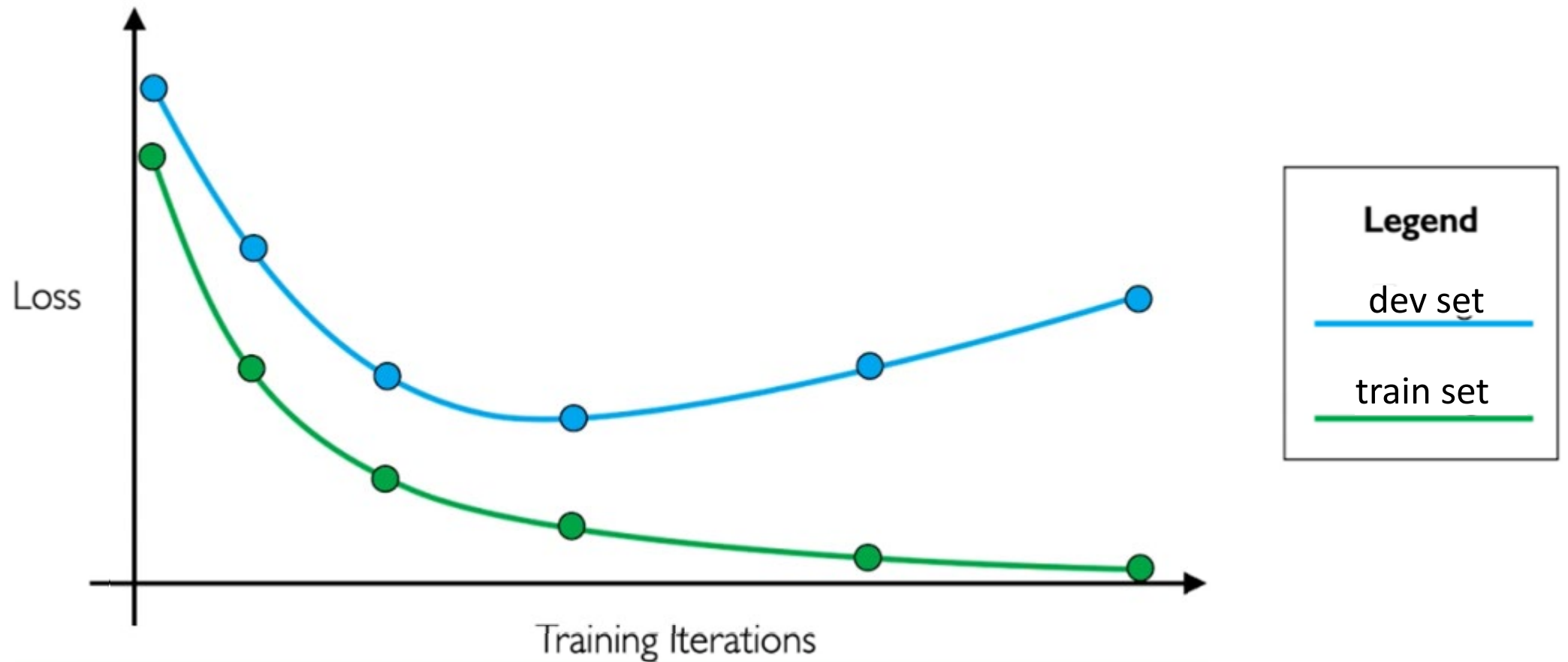**Legend**

dev set ──────

train set ──────

Loss

Training Iterations

# Early Stopping

- Stop training before we have a chance to overfit

# Early Stopping

Orthogonalization

- Optimize cost function $J$
  - gradient
- Not overfit
  - Regularization

$J(\omega, b)$

- Stop training before we have a chance to overfit

$L_2 \rightarrow \lambda$



**Legend**

dev set

train set

$\omega = 0$   mid-size $\|\omega\|_F^2$   large $\omega$

Under-fitting | Over-fitting

Stop training here!

Loss

Training Iterations

# Setting up Your Optimization Problem: Normalizing inputs

# Normalizing inputs

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



subtract mean:

$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$

$$x := x - \mu$$

use $\mu$, $\sigma^2$ to normalize test data.
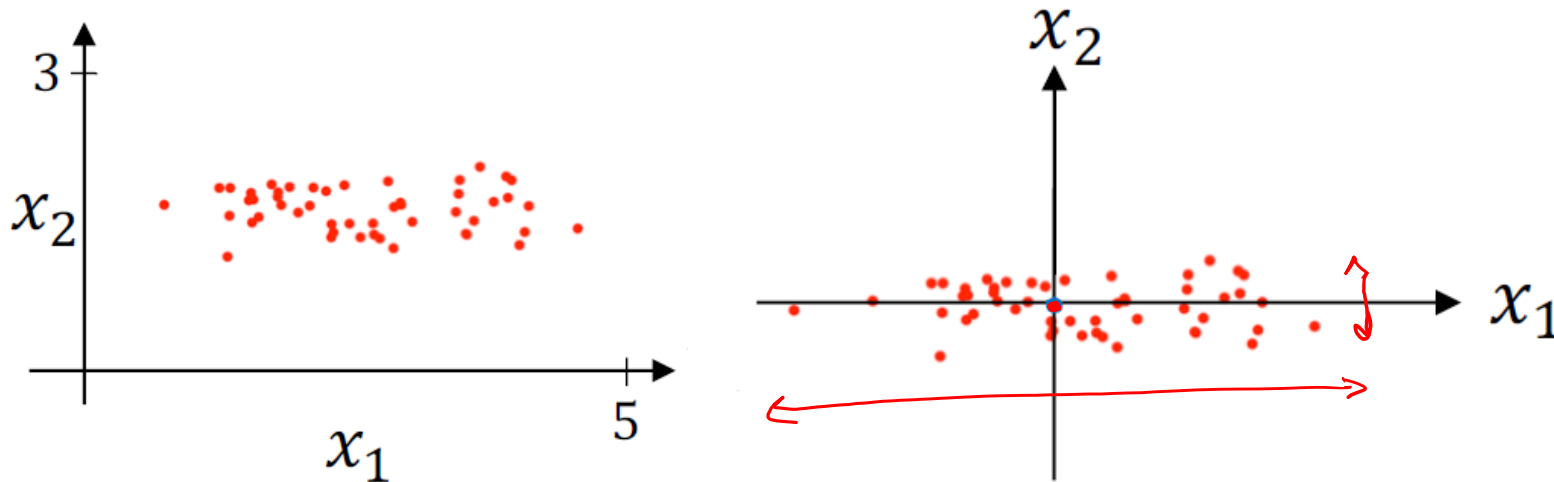
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} **2$$

↑ element-wise

$$x /= \sigma^2$$

↑ normalize variance

# Why normalize inputs?

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

$w \qquad X_1 : 1 , \ldots , 1000$ ←

Unnormalized:

$b \qquad X_2 : 0 \ldots 1$ ←

Normalized:



$X_1 : 0 \ldots 1$

$X_2 : -1 \ldots 1$

$X_3 : 1 \ldots 2$

# Setting up Your Optimization Problem: Vanishing/Exploding Gradients

# Vanishing/Exploding Gradients
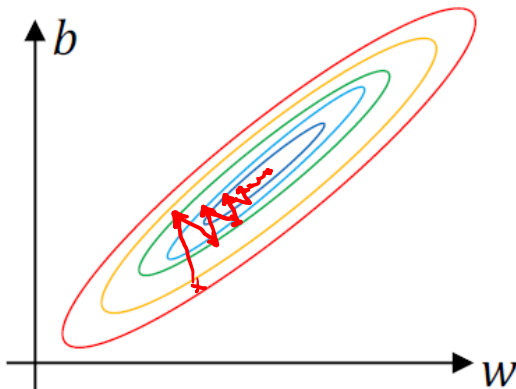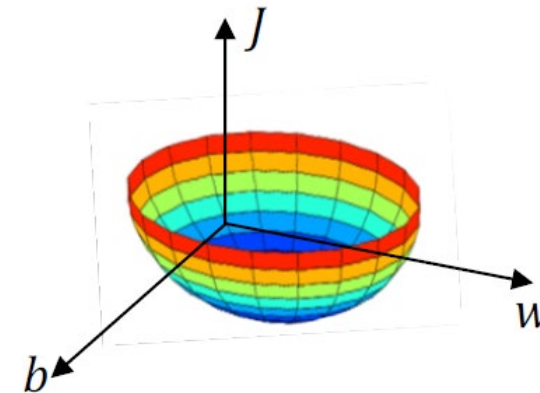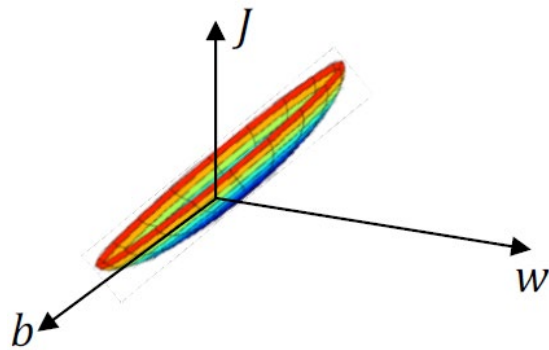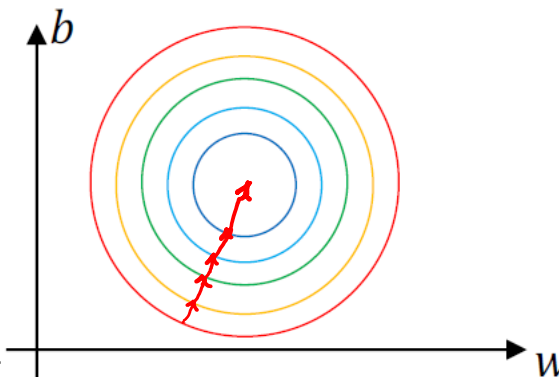
$L = 150$



$x_1$

$x_2$

$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$

$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$

$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$ ... $y_L$

$\hat{y}$

$\omega^{[1]}$ $\omega^{[2]}$ $\omega^{[3]}$

$\omega^{[L]}$

$\frac{L}{1.5}$

$0.5^L$

$$g(z) = z \qquad\qquad b^{[l]} = 0$$

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

$$z^{[1]} = \omega^{[1]} X$$

$$a^{[1]} = g(z^{[1]}) = z^{[1]}$$

$$a^{[2]} = g(z^{[1]}) = g(\omega^{[2]} a^{[1]})$$

$$\omega^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \qquad 1 \le l < L$$

$$\hat{y} = \omega^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} X$$

# Single neuron example

$x_1$

$x_2$

$x_3$

$x_4$

$\hat{y}$

$a = g(z)$

$z = \boxed{\omega_1} x_1 + \boxed{\omega_2} x_2 + \cdots + \boxed{\omega_n} x_n + \cancel{b}$

Large $n \longrightarrow$ smaller $\omega_i$

$var(\omega_i) \approx \dfrac{1}{n}$

$\omega^{[\ell]} = np.random.rand(shape) * np.sqrt\left(\dfrac{2}{n^{[\ell-1]}}\right) \leftarrow$

Relu

tanh

$$\sqrt{\dfrac{1}{n^{[\ell-1]}}} \quad \swarrow$$

Xavier initialization

$$\sqrt{\dfrac{2}{n^{[\ell-1]} + n^{[\ell]}}} \swarrow$$

# Setting up Your Optimization Problem: Numerical approximation of gradients

# Checking your derivative computation

$f(\theta) = \theta^3$



$f(\theta+\varepsilon)$

$f(\theta)$

$f(\theta-\varepsilon)$

$f(\theta+\varepsilon) - f(\theta-\varepsilon)$

$\theta-\varepsilon$    $\theta$    $\theta+\varepsilon$

$0.99$    $1$    $1.01$

$\varepsilon = 0.01$

$$\frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon}$$
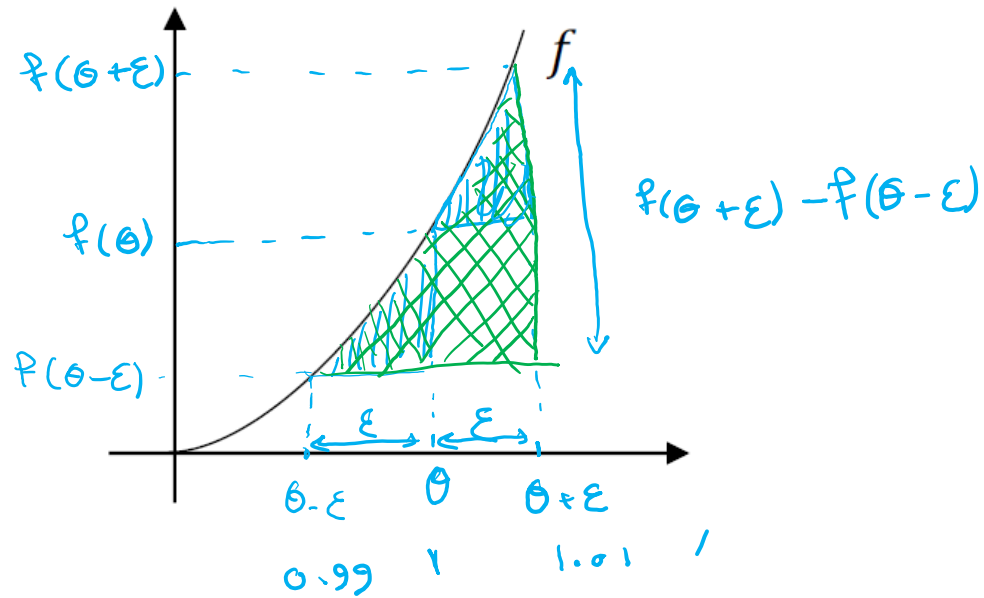
$$\frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon} \simeq g(\theta)$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

appox error :  $0.0001$

$0.1 \rightarrow 0.0001$       error: $0.03$ $\leftarrow$    $g(\theta): 3.0301$

$\frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon} \rightarrow$ error : $O(\varepsilon^2)$        $\frac{f(\theta+\varepsilon) - f(\theta)}{\varepsilon}$     error : $O(\varepsilon)$   $0.01$

# Setting up Your Optimization Problem: Gradient Checking

# Gradient check for a neural network

- Take $w^{[1]}$, $b^{[1]}$ ..., $w^{[L]}$, $b^{[L]}$ and reshape into a big vector $\theta$.

Concatinate

$$J(w^{[1]}, b^{[1]}, \ldots w^{[L]}, b^{[L]}) = J(\theta)$$

- Take $dw^{[1]}$, $db^{[1]}$ ..., $dw^{[L]}$, $db^{[L]}$ and reshape into a big vector $d\theta$.

Concatinate

Is $d\theta$ graidient of $J(\theta)$ ?

# Gradient checking (Grad check)

$$J(\theta) = J(\theta_1, \theta_2, \theta_3 \ldots \theta_i)$$

for each i :

$$d\theta_{approximate}[i] = \frac{J(\theta_1, \theta_2 \ldots \theta_{i+\varepsilon} \ldots) - J(\theta_1, \theta_2 \ldots \theta_{i-\varepsilon} \ldots)}{2\varepsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i}$$

$$d\theta_{approximate} \overset{?}{\approx} d\theta$$

check : $\dfrac{\| d\theta_{app\ldots} - d\theta \|_2}{\| d\theta_{appr} \|_2 + \| d\theta \|_2}$

$$\approx \boxed{\begin{array}{l} 10^{-7} - \text{great !} \\ \hline 10^{-5} \\ 10^{-3} - \text{worry} \end{array}}$$

# Setting up Your Optimization Problem: Gradient Checking Implementation Notes

# Gradient Checking Implementation Notes

- Don't use in training – only to debug

- If algorithm fails grad check, look at components to try to identify bug.

- Remember regularization.

- Doesn't work with dropout.

- Run at random initialization; perhaps again after some training.

$$d\theta_{approx}[i]$$

$$d\theta[i]$$

$$db^{[\ell]}$$

$$J(\theta) = \frac{1}{m} \sum \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{\ell} \|w^{[\ell]}\|_F^2$$

$$\text{keep-prob} = 1.0 \quad d\theta \quad J$$

$$w, b = 0$$

# Core Foundation Review

- Other Regularization Techniques

    - Data Augmentation

    - Early Stopping

- Setting up Your Optimization Problem

    - Normalizing Inputs

    - Vanishing/Exploding Gradients

    - Numerical approximation of gradients

    - Gradient Checking

    - Gradient Checking Implementation Notes