

# Feature engineering

$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1}_{\text{frontage}} \underbrace{x_1}_{\text{frontage}} + \underbrace{w_2}_{\text{depth}} \underbrace{x_2}_{\text{depth}} + b$$

$$\text{area} = \text{frontage} \times \text{depth}$$

$$x_3 = x_1 x_2$$

new feature

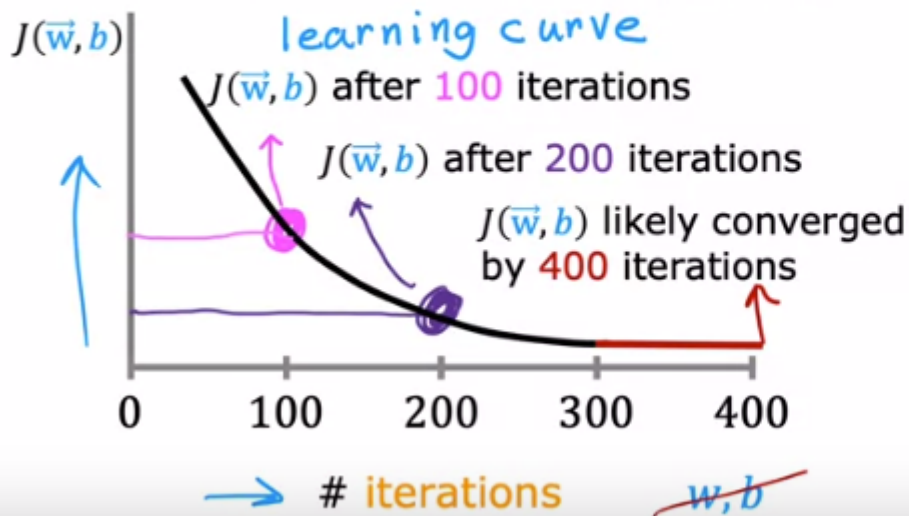
$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1}_{\text{frontage}} x_1 + \underbrace{w_2}_{\text{depth}} x_2 + \underbrace{w_3}_{\text{area}} x_3 + b$$



Feature engineering:  
Using **intuition** to design  
**new features**, by  
transforming or combining  
original features.

# Make sure gradient descent is working correctly

objective:  $\min_{\vec{w}, b} J(\vec{w}, b)$   $J(\vec{w}, b)$  should **decrease** after every iteration



# iterations needed varies 30 1,000 100,000

Automatic convergence test

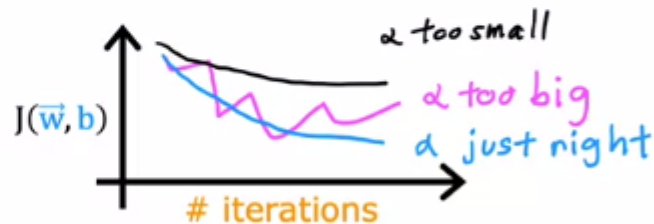
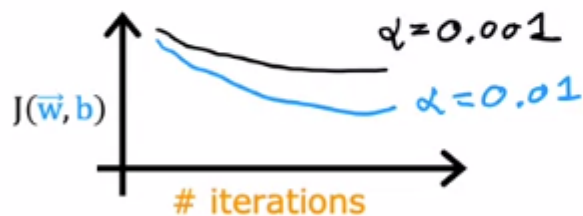
Let  $\epsilon$  "epsilon" be  $10^{-3}$ .  
*0.001*

If  $J(\vec{w}, b)$  decreases by  $\leq \epsilon$  in one iteration,  
declare **convergence**.

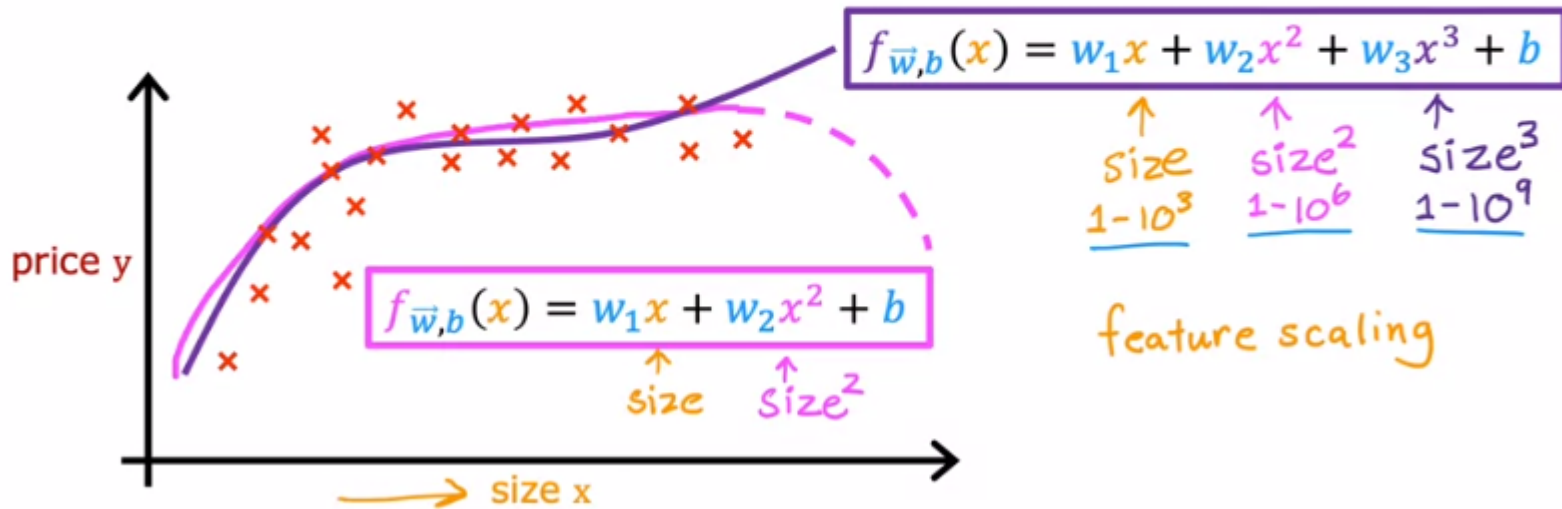
(found parameters  $\vec{w}, b$  to get close to global minimum)

Values of  $\alpha$  to try:

... 0.001 0.003 0.01 0.03 0.1 0.3 1 ...  
           $\nearrow$            $\nearrow$            $\nearrow$            $\nearrow$            $\nearrow$   
           $3\times$            $\approx 3\times$            $3\times$            $\approx 3\times$            $3\times$            $\approx 3\times$



# Polynomial regression



# Feature scaling

aim for about  $-1 \leq x_j \leq 1$  for each feature  $x_j$

$-3 \leq x_j \leq 3$   
 $-0.3 \leq x_j \leq 0.3$  } acceptable ranges

$$0 \leq x_1 \leq 3$$

okay, no rescaling

$$-2 \leq x_2 \leq 0.5$$

okay, no rescaling

$$-100 \leq x_3 \leq 100$$

too large  $\rightarrow$  rescale

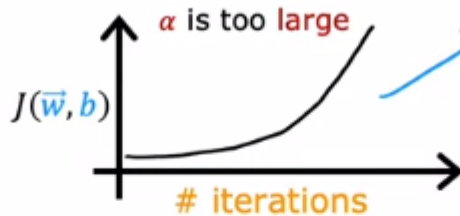
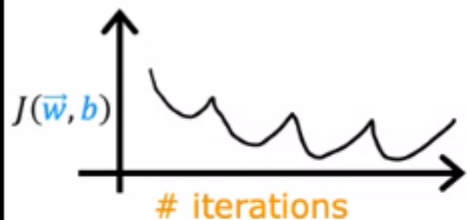
$$-0.001 \leq x_4 \leq 0.001$$

too small  $\rightarrow$  rescale

$$98.6 \leq x_5 \leq 105$$

too large  $\rightarrow$  rescale

## Identify problem with gradient descent



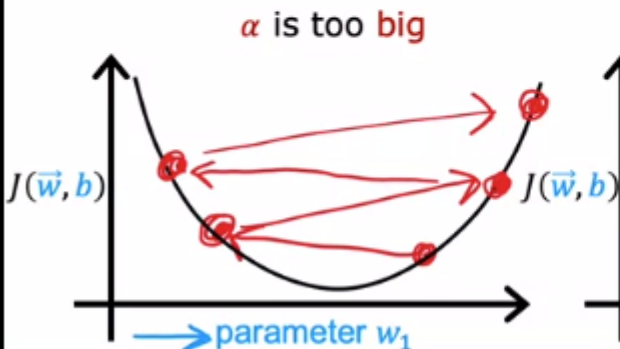
or learning rate is too large

$$w_1 = w_1 + \alpha d_1 \quad \text{!!}$$

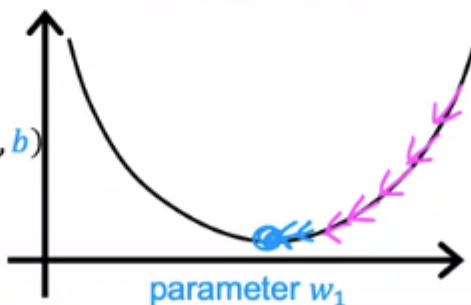
use a minus sign

$$w_1 = w_1 - \alpha d_1 \quad \text{!!}$$

## Adjust learning rate



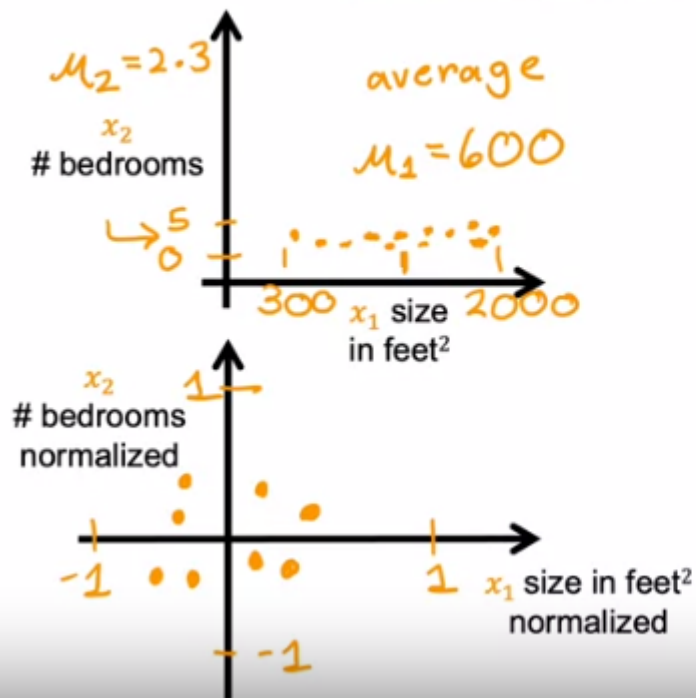
Use smaller  $\alpha$



With a small enough  $\alpha$ ,  $J(\bar{w}, b)$  should **decrease** on every iteration

If  $\alpha$  is too small, gradient descent takes a lot more iterations to **converge**

# Mean normalization



$$300 \leq x_1 \leq 2000$$

$$x_1 = \frac{x_1 - \mu_1}{2000 - 300}$$

max-min

$$-0.18 \leq x_1 \leq 0.82$$

$$0 \leq x_2 \leq 5$$

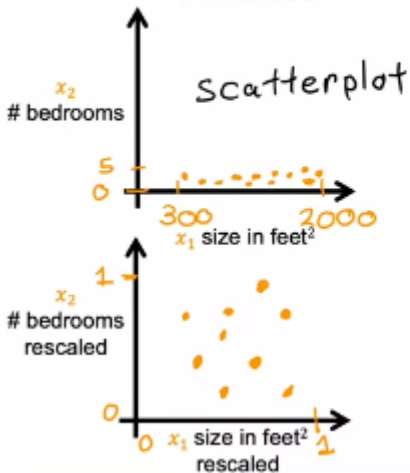
$$x_2 = \frac{x_2 - \mu_2}{5 - 0}$$

max-min

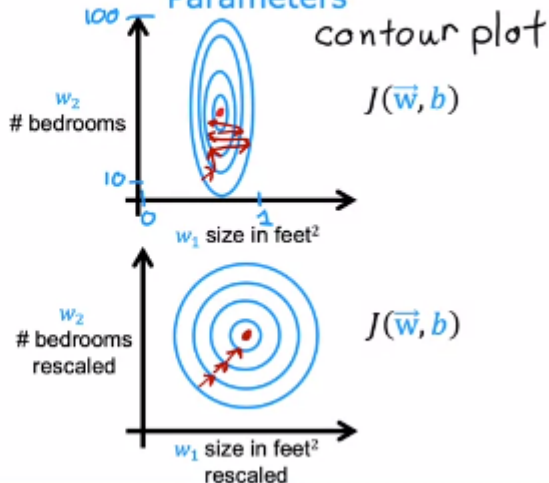
$$-0.46 \leq x_2 \leq 0.54$$

# Feature size and gradient descent

Features

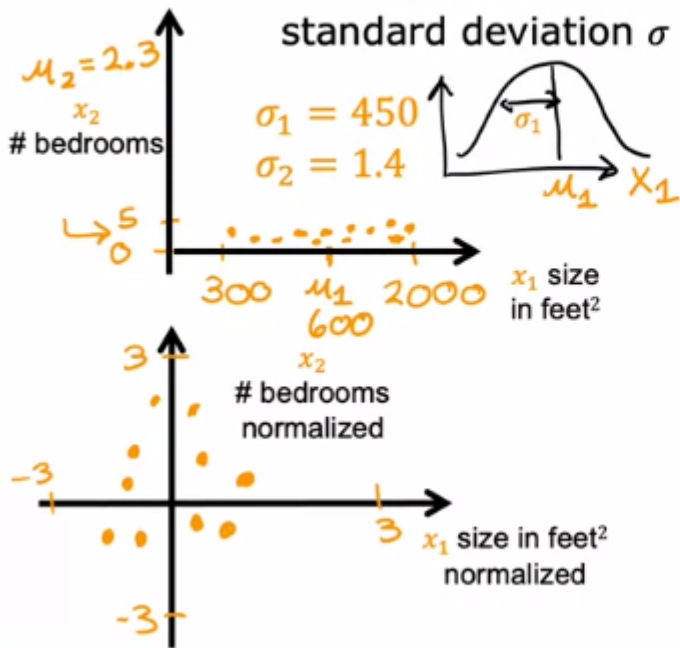


Parameters





# Z-score normalization



$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

$$x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

$$-0.67 \leq x_1 \leq 3.1 \quad -1.6 \leq x_2 \leq 1.9$$