# Multiple features

| $x_1$ size ($ft^2$) | # of bed $x_2$ | # $y$ bath $x_3$ | age $x_4$ | Price $y$ ($\$1000's$) |
|---|---|---|---|---|
| — | | | | — |
| — | | | | — |
| — | | | | — |
| — | | | | — |
| — | | | | — |
| — | | | | — |

1) One feature: $f_{w,b} = wx+b$

2) Multiple features:

$x_j = j^{th}$ feature

$\vec{x}^{(i)} =$ features of $i^{th}$ training sample $\Rightarrow$ e.g. $\vec{x}^{(3)} = \begin{bmatrix} 2000 & 3 & 2 & 30 & 5^{00} \end{bmatrix}$

$x_j^{(i)} =$ value of $j^{th}$ feature in $i^{th}$ training sample $\Rightarrow x_4^{(3)} = 30$

$f_{x,b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$

e.g. : $\quad f_{x,b}(\vec{x}) = 0.1 \, x_1 + 4 x_2 + 3 x_3 + (-2)x_4 + 80$

$\qquad\qquad\qquad$ size $\quad$ #bed $\quad$ #bath $\quad$ age $\quad$ base price

Mac General:

$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$

$\vec{w} = \begin{bmatrix} w_1 & w_2 & \dots & w_n \end{bmatrix} \leftarrow$ model parameters

$b =$ number

$\vec{x} = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}$

more compact in matrix/vector (linear algebra) expression:

$f_{\vec{w},b} = \vec{w} \cdot \vec{x} + b$

$\quad$ dot product : $\vec{w} \cdot \vec{x} = \sum_j w_i x_i$

$\qquad = \sum_i w_i x_i + b$

"multiple linear regression" not "multivariate regression"

Vectorization: (code)

$x[1] \quad x[2] \quad x[3]$

*) w/o. vectorization / loop : $f = w[0] + w[1] + w[2] + w[3] + b$

$\Rightarrow$ cumbersome

2) loop : $f = b$

for j in range (0,n): $\Leftarrow$ n steps to compute

$f = f + w[j] * x[j]$

3) Vectorized : $f = np.dot(w,x) + b \Leftarrow 1$ step! (parallelized)

1) then, Gradient descent we have

$$\vec{w} = \vec{w} - \alpha \vec{d}$$

$\underset{1 \text{ op}}{\uparrow} \quad \underset{1 \text{ op}}{\uparrow}$

// Gradient Descent for multiple Linear Regression:

Vector notation:

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

$$J(\vec{w},b) = \frac{1}{2m} \sum_{i=1}^{i=m} \left[ (\vec{w} \cdot \vec{x}^{(i)} + b) - y^{(i)} \right]^2$$

update rule:

$\Rightarrow \quad w_j = w_j - \alpha \underbrace{\frac{d}{dw_j} J(\vec{w},b)}$

$b = b - \alpha \frac{d}{db} J(\vec{w},b)$

$\Rightarrow \quad \vec{w} = \vec{w} - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{w,b}(\vec{x}^{(i)} - y^{(i)}) \right) \vec{x}^{(i)}$ $\Big\}$ simultaneous update

$b = b - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right)$

## Alternative to gradient descent ( for Linear Regression Only)

- Solve for $w, b$ w/o interations
  ( symbolically ) $\Longrightarrow$ $\vec{\theta} = (x^T \cdot x)^{-1} x^T \cdot y$

- Slow when large # of features

- May be used in some ML packages

- GD is recommended.

## Feature Scaling

Example : $\hat{y} = w_1 x_1 + w_2 x_2 + b$

price ↗    size ↑    # bed ↑

$\sim 200 = 500.$ range
$\Rightarrow$ big

1 - 5 range
$\Rightarrow$ small

Example data point

$(2000, 5, 500)$

size ↑   #br ↑   price ↑

If    $w_1 = 50$    $w_2 = 0.1$    $b = 50$
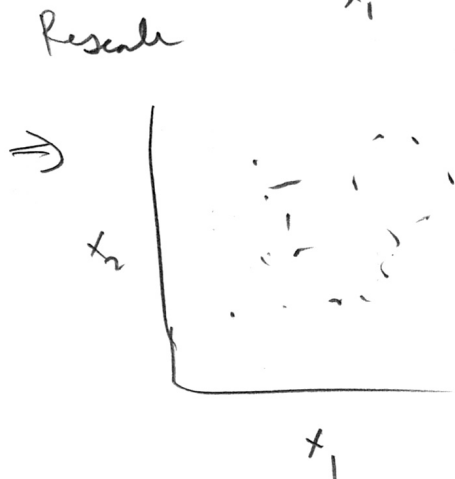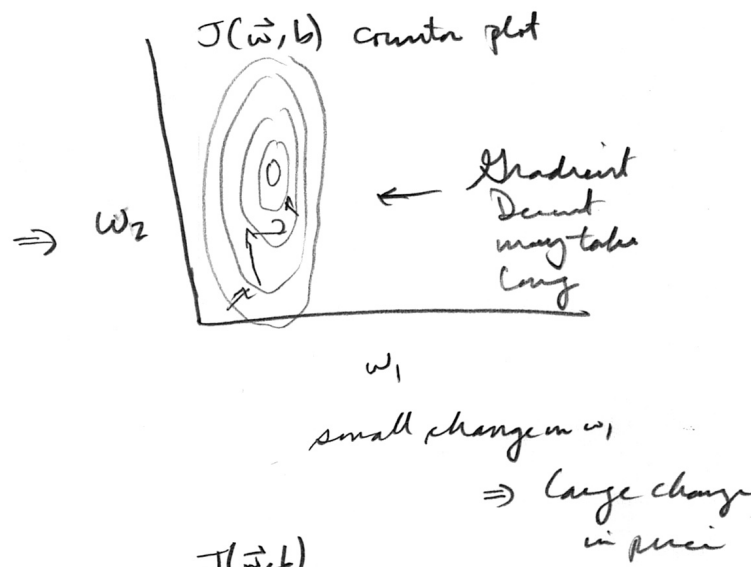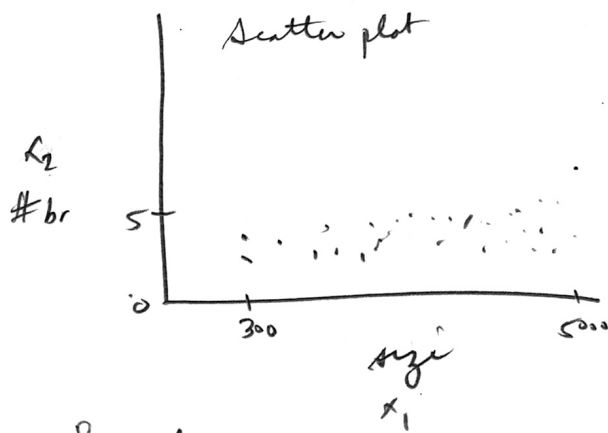
$\hat{y} = \underbrace{50 \cdot 2000}_{\$100,000 k} + 0.1 \cdot 5 + 50$

$\$0.5k + \$50k = \$1,009,050,5 \Rightarrow$ too big

If    $w = 0.1$    $w_2 = 50$    $b = 50$

$\hat{y} = 0.1 \cdot 2000k + 50 \cdot 5 + 50$

      200k      250 k     $50k = \$500k$

$\Rightarrow$ small $w$ for large $x$ range

Scatter plot

$x_2$
#br   5

0   300          5000

size
$x_1$

$J(\vec{w}, b)$ contour plot

$\Rightarrow w_2$

← Gradient
Descent
may take
long

$w_1$

small change in $w_1$
$\Rightarrow$ large change
in price

Rescale

$\Rightarrow$

$x_2$

$x_1$

$J(\vec{w}, b)$

$\Rightarrow w_2$

$w_1$

=

• Feature Scaling    (max normalization)

$300 \leq x_1 \leq 5000$          $1 \leq x_2 \leq 5$

$\Rightarrow x_{1, scaled} = \dfrac{x_1}{5000}$     $x_{2, scaled} = \dfrac{x_2}{5}$

$\Rightarrow 0.15 \leq x_{1, scaled} \leq 1$     $.2 \leq x_{2, scaled} \leq 1$

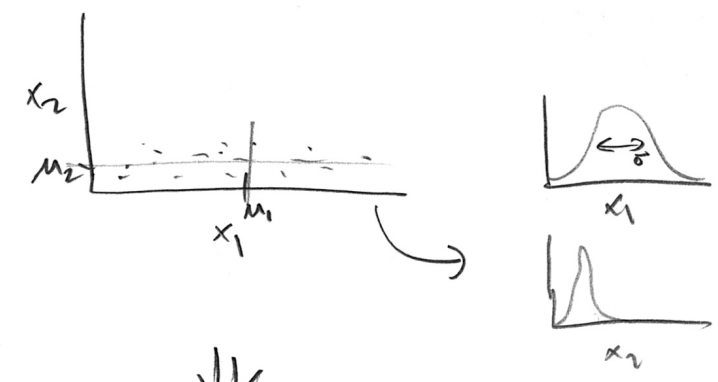• Mean normalization          compute mean $\mu_1, \mu_2$ of $x_1, x_2$

$$x_{1,s} = \frac{x_1 - \mu_1}{5000 - 300} \qquad x_2 = \frac{x_2 - \mu_2}{5 - 1}$$

↑           ↑
max        min

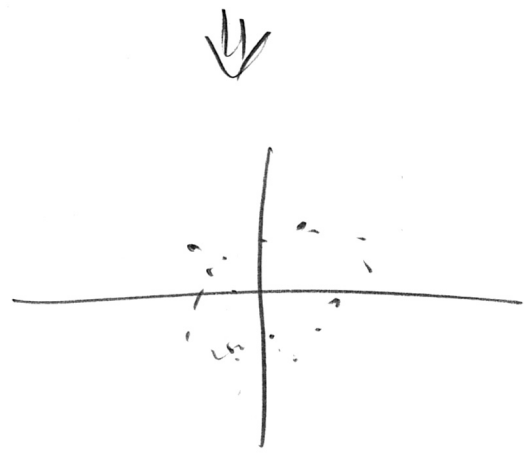$\Rightarrow -0.18 \leq x_1 \leq 0.82 \qquad -0.46 \leq x_2 \leq 0.54$

# Z-Scale Normalization



compute the standard deviation $\sigma_1, \sigma_2$ of $x_1, x_2$.

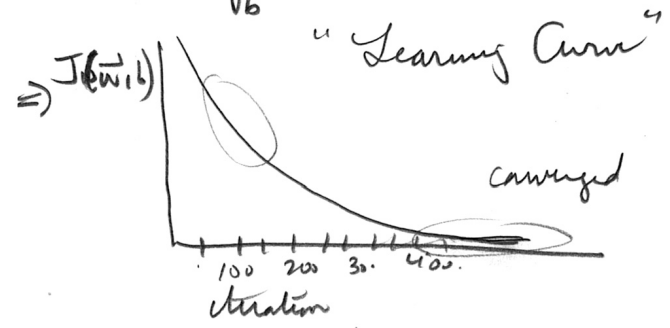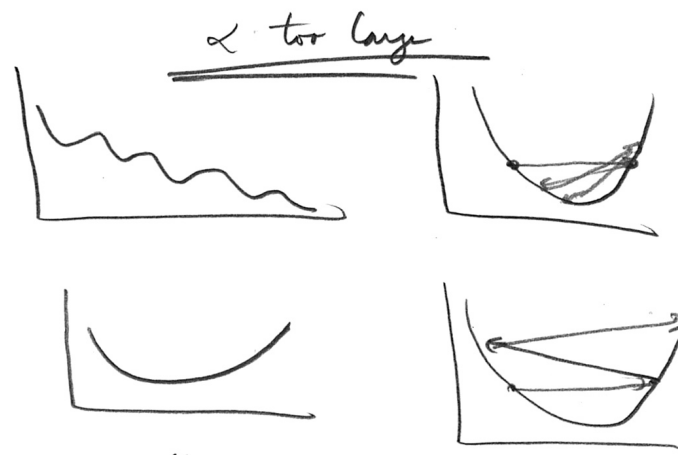$$\Rightarrow x_{1,s} = \frac{x_1 - \mu_1}{\sigma_1}$$

$$x_{2,s} = \frac{x_2 - \mu_2}{\sigma_2}$$

## Gradient Descent Convergence:

$$\omega_{j+1} = \omega_j - \alpha \frac{d}{d\omega_j} J(\vec{w}, b)$$

$$b_{j+1} = b - \alpha \frac{d}{db} J(\vec{w}, b)$$

$\alpha$ too large



$\Rightarrow$ try different values of $\alpha$
$10^{-6} - 10^{-1}$

"Learning Curve"

$\Rightarrow J(\vec{w}, b)$

converged $\Rightarrow$ by eye!

100  200  3..  4..
iteration

Automatic (not perfect)

Compute $\epsilon = \omega_j - \omega_{j-1}$

Converge when $\epsilon$ small e.g. 0.001
("declare convergence")

# Feature Engineering

Example    $x_1$    length    of property
           $x_2$    width

$$\Rightarrow f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x_2 + b$$

But price more likely go as area $= x_1 \cdot x_2 \leftarrow$ new
feature
"$x_3$"

$$\Rightarrow f_{\vec{w},b}(\vec{x}) = w_1 x_1 * w_2 x_2 + w_3 x_3 + b$$

$\uparrow$
$x_1 x_2$

Feature Engineering:   use intuition to
                       design new features
                       (transform / combine other
                               features)

# Polynomial Regression



$\leftarrow f_{w,b}(\vec{x}) = w_1 x + w_2 x^2 + w_3 x^3 + b$

$\nwarrow f_{\vec{w},b}(\vec{x}) = w_1 x_1 + b$

$f(\vec{w},b(\vec{x}) = w_1 x_1 + w_2 x_1^2 + b$

$f_{\vec{w},b}(x) = w_1 x + w_2 \sqrt{x} + b$