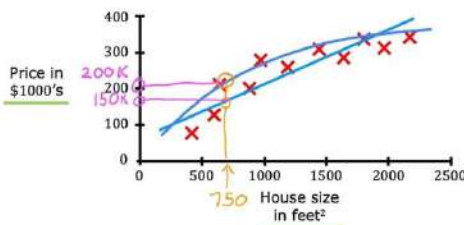# Supervised Machine Learning

## Regression

predict a number from many possible outputs

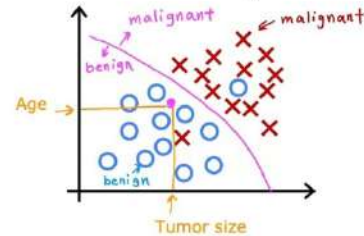## Classification

small number of outputs

tumor → benign
tumor → malignant



Regression: Housing price prediction
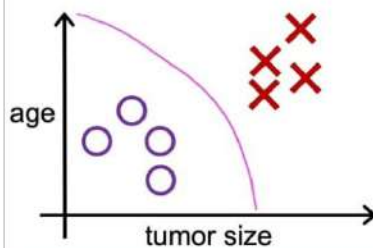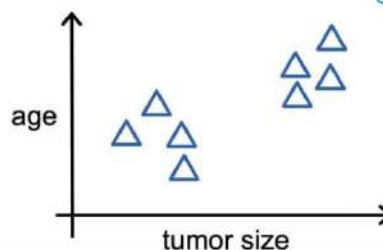


Two or more inputs

# Supervised vs Unsupervised



Supervised learning
Learn from data labeled with the "right answers"

Unsupervised learning
Find something interesting in unlabeled data.

# Unsupervised Machine Learning

**Clustering**

take data without labels and automatically group them to clusters

**Anomaly Detection**

Find unusual data points

**Dimensionality Reduction**

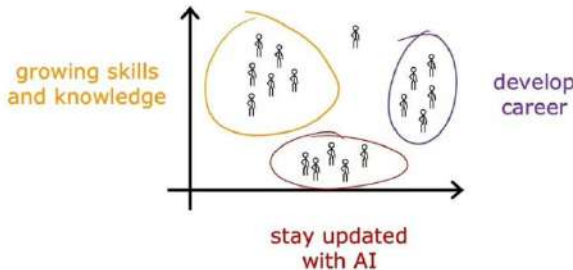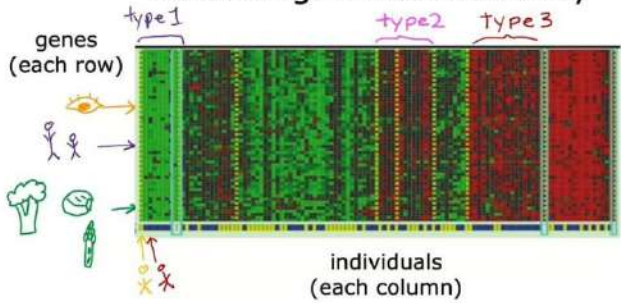Compress data using fewer numbers



Clustering: Google news



Clustering: Grouping customers

growing skills and knowledge

develop career

stay updated with AI



Clustering: DNA microarray

type1 type2 type3

genes (each row)

individuals (each column)

# Linear Regression

## House sizes and prices



linear regression

price in $1000's — **Numbers**

220 K

1.5
−33.2

size in feet²

1250

categories
cat } 2
dog

disease ✚ 10

**Regression model**
Predicts numbers
Infinitely many possible outputs

**Supervised learning model**
Data has "right answers"

**Classification model**
Predicts categories
Small number of possible outputs

---

## Terminology

Training set: Data used to train the model

| | $x$ size in feet² | $y$ price in $1000's |
|---|---|---|
| (1) | 2104 | 400 |
| (2) | 1416 | 232 |
| (3) | 1534 | 315 |
| (4) | 852 | 178 |
| ... | ... | ... |
| (47) | 3210 | 870 |

$m = 47$

$x^{(1)} = 2104 \qquad y^{(1)} = 400$

$(x^{(1)}, y^{(1)}) = (2104, 400)$

$x^{(2)} = 1416 \qquad x^{(2)} \neq x^2$ not exponent

**Notation:**

$x$ = "input" variable feature

$y$ = "output" variable "target" variable

$m$ = number of training examples

$(x, y)$ = single training example

$(x^{(i)}, y^{(i)})$

$(x^{(i)}, y^{(i)})$ = $i^{th}$ training example

index        (1st, 2nd, 3rd ...)

---



training set  → features / targets

learning algorithm

$x$ → $f$ → $\hat{y}$   "y-hat"  $\hat{y}$

~~hypothesis~~  function

feature   model   prediction (estimated $y$)

target

size → $f$ → price (estimated)

## How to represent $f$?

$f_{w,b}(x) = wx + b$

$f(x)$

$f_{w,b}(x) = wx + b$

$f(x) = wx + b$

linear

single feature $x$

Linear regression with **one** variable.
size

**Univariate** linear regression.
one variable

# Cost function

## Cost function: Squared error cost function



$$J(w,b) = \frac{1}{2m}\sum_{i=1}^{m}\left(\hat{y}^{(i)} - y^{(i)}\right)^2 \quad \text{error}$$

$m$ = number of training examples

$$J(w,b) = \frac{1}{2m}\sum_{i=1}^{m}\left(f_{w,b}(x^{(i)}) - y^{(i)}\right)^2$$
↑ intuition

$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$ ←

$f_{w,b}(x^{(i)}) = wx^{(i)} + b$

Find $w, b$:
$\hat{y}^{(i)}$ is close to $y^{(i)}$ for all $(x^{(i)}, y^{(i)})$.

---

### model:
$$f_{w,b}(x) = wx + b$$

### parameters:
$w, b$

### cost function:
$$J(w,b) = \frac{1}{2m}\sum_{i=1}^{m}(f_{w,b}(x^{(i)}) - y^{(i)})^2$$

### goal:
$$\min_{w,b} J(w,b)$$

### simplified
$$f_w(x) = wx \qquad b = \emptyset$$

$w$

$$J(w) = \frac{1}{2m}\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})^2$$

$$\min_{w} J(w)$$
$wx^{(i)}$

---

$f_w(x)$
(for fixed $w$, function of $x$)
input

$J(w)$
(function of $w$)
parameter

$w = 1$
$f(1) = 1$
$f(x) = y$
$y = 1$    $x = 1$

$J(1) = 0$

$w = 1$
$$J(w) = \frac{1}{2m}\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m}\sum_{i=1}^{m}(wx^{(i)} - y^{(i)})^2 = \frac{1}{2m}(0^2 + 0^2 + 0^2) = 0$$
$wx^{(i)}$        $\emptyset$

---

$f_w(x)$
(function of $x$)
$w = 1$
$f_w$
$f$
$w = 0$
$f(x) = -0.5x$

5.25

$J(w)$
(function of $w$)
$w = 1$
$J$
$w$

$$J(0) = \frac{1}{2m}\left(1^2 + 2^2 + 3^2\right) = \frac{1}{6}[14] \approx 2.3$$

# Visualize cost function



Model
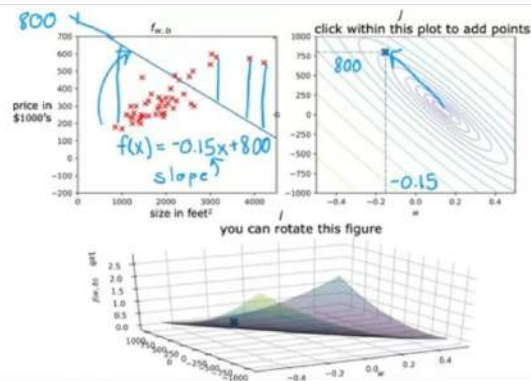
$$f_{w,b}(x) = wx + b$$

Parameters

$$w, b$$

Cost Function

$$J(w,b) = \frac{1}{2m}\sum_{i=1}^{m}(f_{w,b}(x^{(i)}) - y^{(i)})^2$$

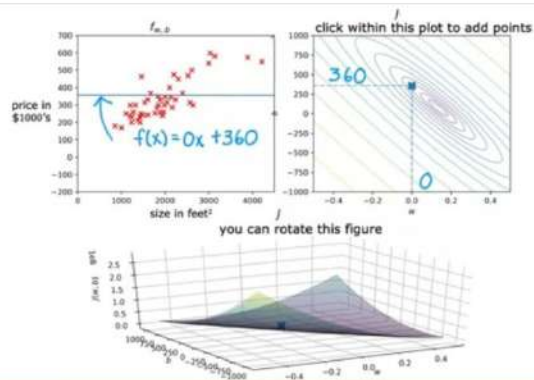Objective

$$\underset{w,b}{\text{minimize}}\, J(w,b)$$

800

click within this plot to add points

$f(x) = -0.15x + 800$

slope

-0.15

you can rotate this figure

click within this plot to add points

360

$f(x) = 0x + 360$

0

you can rotate this figure

click within this plot to add points

error
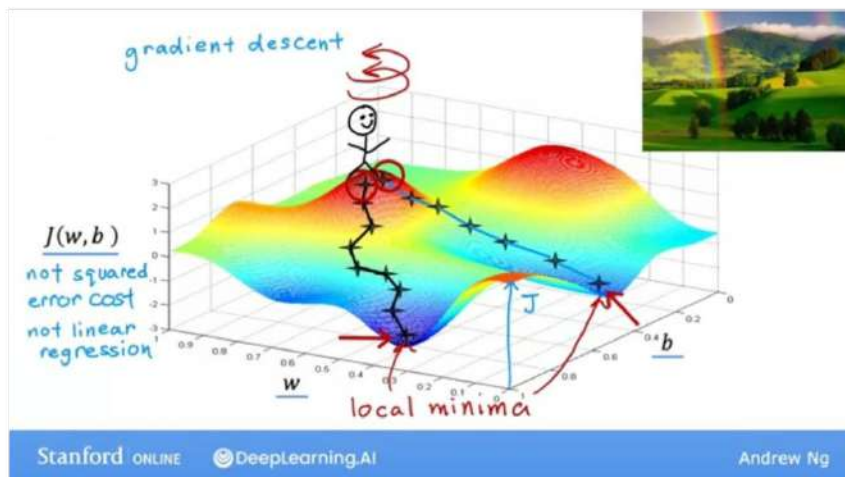
71

$f(x) = 0.13x + 71$

0.13

you can rotate this figure

# Gradient Descent

---

Have some function $J(w,b)$ *for linear regression or any function*

Want $\min_{w,b} J(w,b)$
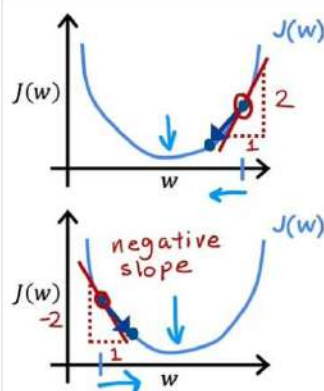
$$\min_{w_1,\ldots,w_n,b} J(w_1, w_2, \ldots, w_n, b)$$

Outline:

Start with some $w,b$ (set $w=0$, $b=0$)

Keep changing $w,b$ to reduce $J(w,b)$    *J not always*

Until we settle at or near a minimum

*may have >1 minimum*

---

*gradient descent*

$J(w,b)$
*not squared error cost*
*not linear regression*

*w*     *b*

*J*

*local minima*

---

## Gradient descent algorithm

repeat until convergence {

*learning rate*

$$w = w - \alpha \frac{\partial}{\partial w} J(w,b) \quad \text{derivative}$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w,b)$$

}

$J(w)$

$$w = w - \alpha \frac{\partial}{\partial w} J(w)$$

$$\min_w J(w)$$

---

$J(w)$

$J(w)$

2

1

$w$

*negative slope*   $J(w)$

$J(w)$

-2

1

$w$

$$w = w - \alpha \frac{d}{dw} J(w)$$

$> 0$

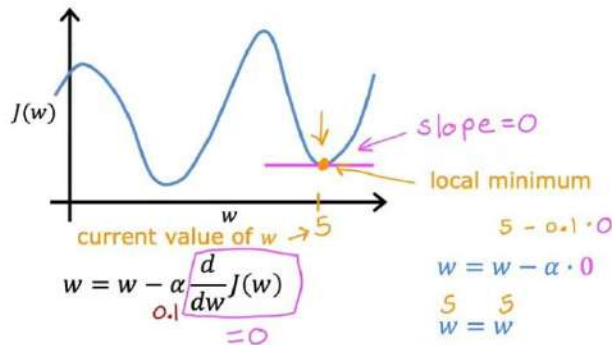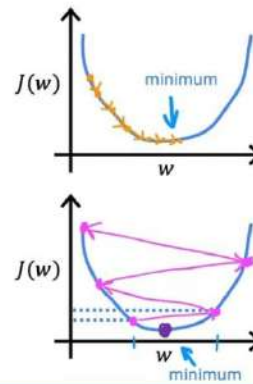$$w = w - \alpha \cdot (\text{positive number})$$

$$\frac{d}{dw} J(w) < 0$$

$$w = w - \alpha \cdot (\text{negative number})$$

$$w = w - \alpha \frac{d}{dw} J(w)$$

If $\alpha$ is too small...
Gradient descent may be slow.

If $\alpha$ is too large...

Gradient descent may:
- Overshoot, never reach minimum
- F

slope = 0

local minimum

current value of $w$ → 5

$$w = w - \alpha \frac{d}{dw} J(w)$$
0.1              = 0

$$5 - 0.1 \cdot 0$$
$$w = w - \alpha \cdot 0$$
$$\underset{5}{w} = \underset{5}{w}$$

# Gradient Disecent for Linear Regression

Linear regression model          Cost function
$$f_{w,b}(x) = wx + b \qquad J(w,b) = \frac{1}{2m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)})^2$$
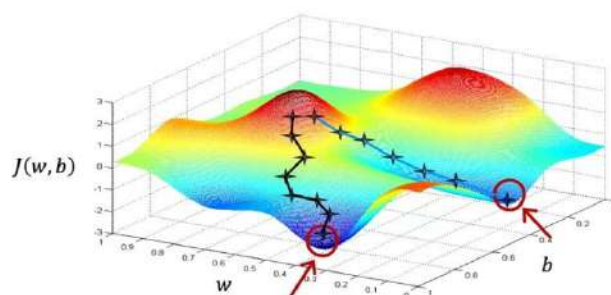
Gradient descent algorithm
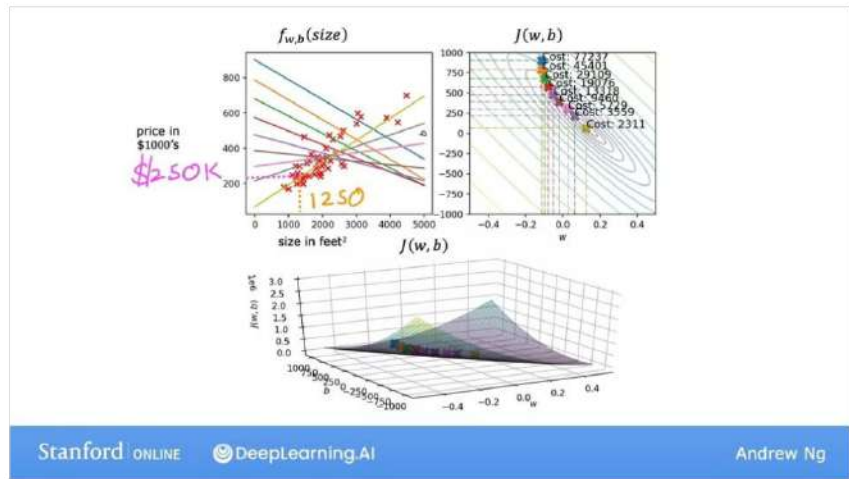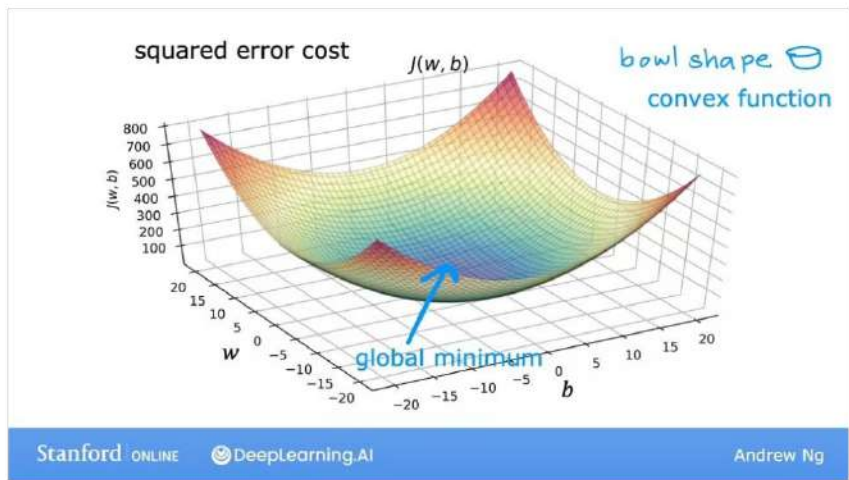repeat until convergence {
$$w = w - \alpha \frac{\partial}{\partial w} J(w,b) \rightarrow \frac{1}{m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)}$$
$$b = b - \alpha \frac{\partial}{\partial b} J(w,b) \rightsquigarrow \frac{1}{m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)})$$
}

More than one local minimum



$J(w,b)$

$w$                          $b$

# Week 2

## Multiple Linear Regression

### Multiple features (variables)

| Size in feet² | Number of bedrooms | Number of floors | Age of home in years | Price ($) in $1000's |
|---|---|---|---|---|
| $X_1$ | $X_2$ | $X_3$ | $X_4$ | |
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

$i = 2$ (row), $j = 1 \dots 4$, $n = 4$

$x_j = j^{th}$ feature
$n$ = number of features
$\vec{x}^{(i)}$ = features of $i^{th}$ training example
$x_j^{(i)}$ = value of feature $j$ in $i^{th}$ training example

$$\vec{x}^{(2)} = \begin{bmatrix} 1416 & 3 & 2 & 40 \end{bmatrix}$$

$$x_3^{(2)} = 2$$

### Model:

Previously: $f_{w,b}(x) = wx + b$

$$f_{w,b}(X) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$$

example

$$f_{w,b}(x) = 0.1 x_1 + 4 x_2 + 10 x_3 + -2 x_4 + 80$$

↑ size  ↑ #bedrooms  ↑ #floors  ↑ years  ↑ base price

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

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

$$\vec{w} = \begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_n \end{bmatrix}$$ parameters of the model

$b$ is a number

Vector $\vec{x} = \begin{bmatrix} x_1 & x_2 & x_3 & \dots & x_n \end{bmatrix}$

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1 x_1 + w_2 x_2 + w_3 x_3 + \cdots + w_n x_n + b$$

↑ dot product    multiple linear regression

(not multivariate regression)

# Multiple Features

## Multiple features (variables)

|  | Size in feet² | Number of bedrooms | Number of floors | Age of home in years | Price ($) in $1000's |
|---|---|---|---|---|---|
|  | $X_1$ | $X_2$ | $X_3$ | $X_4$ |  |
|  | 2104 | 5 | 1 | 45 | 460 |
| $i=2$ | 1416 | 3 | (2) | 40 | 232 |
|  | 1534 | 3 | 2 | 30 | 315 |
|  | 852 | 2 | 1 | 36 | 178 |
|  | ... | ... | ... | ... | ... |

$j = 1 \dots 4$

$n = 4$

$x_j = j^{th}$ feature
$n$ = number of features
$\vec{x}^{(i)}$ = features of $i^{th}$ training example
$x_j^{(i)}$ = value of feature $j$ in $i^{th}$ training example

$\vec{x}^{(2)} = \begin{bmatrix} 1416 & 3 & ②  & 40 \end{bmatrix}$

$x_3^{(2)} = 2$

## Model:

Previously: $f_{w,b}(x) = wx + b$

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

example

$f_{w,b}(x) = 0.1 x_1 + 4 x_2 + 10 x_3 + -2 x_4 + 80$

size   #bedrooms   #floors   years   base price

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

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

$\vec{w} = \begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_n \end{bmatrix}$  parameters of the model

$b$ is a number

vector $\vec{x} = \begin{bmatrix} x_1 & x_2 & x_3 & \dots & x_n \end{bmatrix}$

$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1 x_1 + w_2 x_2 + w_3 x_3 + \cdots + w_n x_n + b$

dot product   multiple linear regression

(not multivariate regression)

# Vectorization

## Parameters and features

$\vec{w} = [w_1 \quad w_2 \quad w_3]$ $\quad n=3$

$b$ is a number

$\vec{x} = [x_1 \quad x_2 \quad x_3]$

linear algebra: count from 1

**NumPy**

$w[\bullet] \quad w[1] \quad w[2]$

```
w = np.array([1.0,2.5,-3.3])
b = 4                x[0]  x[1]  x[2]
x = np.array([10,20,30])
```
code: count from 0

## Without vectorization   $n = 100,000$

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

```
f = w[0] * x[0] +
    w[1] * x[1] +
    w[2] * x[2] + b
```

## Without vectorization

$f_{\vec{w},b}(\vec{x}) = \left( \sum_{j=1}^{n} w_j x_j \right) + b$

$\sum_{j=1}^{n} \rightarrow j=1\dots n$
$\quad 1,2,3$

$range(0,n) \rightarrow j = 0 \dots n-1$

```
f = 0        range(n)
for j in range(0,n):
    f = f + w[j] * x[j]
f = f + b
```

## Vectorization

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

```
f = np.dot(w,x) + b
```

Stanford ONLINE  DeepLearning.AI          Andrew Ng

---

## Without vectorization

```
for j in range(0,16):
    f = f + w[j] * x[j]
```
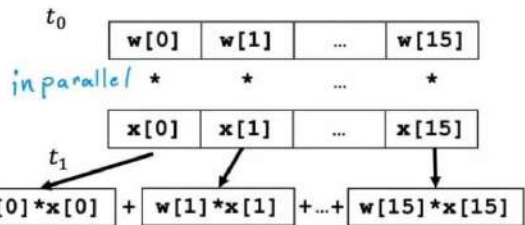
$t_0$
    f + w[0] * x[0]

$t_1$
    f + w[1] * x[1]

...

$t_{15}$
    f + w[15] * x[15]

## Vectorization

```
np.dot(w,x)
```

$t_0$

| w[0] | w[1] | ... | w[15] |

in parallel  *     *    ...    *

| x[0] | x[1] | ... | x[15] |

$t_1$

| w[0]*x[0] | + | w[1]*x[1] | +...+ | w[15]*x[15] |

efficient → scale to large datasets

Stanford ONLINE  DeepLearning.AI          Andrew Ng

---

## Gradient descent

$\vec{w} = (w_1 \quad w_2 \quad \cdots \quad w_{16})$  parameters

derivatives $\vec{d} = (d_1 \quad d_2 \quad \cdots \quad d_{16})$

```
w = np.array([0.5, 1.3, … 3.4])
d = np.array([0.3, 0.2, … 0.4])
```
learning rate $\alpha$

compute $w_j = w_j - 0.1 d_j$ for $j = 1 \dots 16$

### Without vectorization

$w_1 = w_1 - 0.1 d_1$
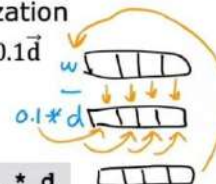$w_2 = w_2 - 0.1 d_2$
$\vdots$
$w_{16} = w_{16} - 0.1 d_{16}$

```
for j in range(0,16):
    w[j] = w[j] - 0.1 * d[j]
```

### With vectorization

$\vec{w} = \vec{w} - 0.1\vec{d}$

$w$

$0.1 * d$

```
w = w - 0.1 * d
```

Stanford ONLINE  DeepLearning.AI          Andrew Ng

# Gradient Descent for multiple linear regression

## Slide 1

| | Previous notation | Vector notation |
|---|---|---|
| Parameters | $w_1, \cdots, w_n$ <br> $b$ | *vector of length n* <br> $\vec{w} = [w_1 \ \cdots \ w_n]$ <br> $b$ still a number |
| Model | $f_{\vec{w},b}(\vec{x}) = w_1 x_1 + \cdots + w_n x_n + b$ | $f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$ <br> *dot product* |
| Cost function | $J(w_1, \cdots, w_n, b)$ | $J(\vec{w}, b)$ |

Gradient descent

repeat {
$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_1, \cdots, w_n, b)$$
$$b = b - \alpha \frac{\partial}{\partial b} J(w_1, \cdots, w_n, b)$$
}

repeat {
$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$$
$$b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$$
}

## Slide 2

### Gradient descent

**One feature**

repeat {
$$w = w - \alpha \frac{1}{m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$
$$\hookrightarrow \frac{\partial}{\partial w} J(w, b)$$
$$b = b - \alpha \frac{1}{m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)})$$

simultaneously update $w, b$
}

**$n$ features ($n \geq 2$)**

repeat {
$j=1$
$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)}) x_1^{(i)}$$
$$\hookrightarrow \frac{\partial}{\partial w_1} J(\vec{w}, b)$$
$\vdots$
$j=n$
$$w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)}$$
$$b = b - \alpha \frac{1}{m} \sum_{i=1}^{m} (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)})$$

simultaneously update
$w_j$ (for $j = 1, \cdots, n$) and $b$
}

## Slide 3

### An alternative to gradient descent

→ Normal equation
- Only for linear regression
- Solve for w, b without iterations

Disadvantages
- Doesn't generalize to other learning algorithms.
- Slow when number of features is large (> 10,000)

What you need to know
- Normal equation method may be used in machine learning libraries that implement linear regression.
- Gradient descent is the recommended method for finding parameters w,b

# Feature Scaling

## Feature and parameter values

$$\widehat{price} = w_1 x_1 + w_2 x_2 + b$$

size ↓   # bedrooms ↓

$x_1$: size (feet²)     $x_2$: # bedrooms
range: $300 - 2,000$    range: $0 - 5$
large                   small

House: $x_1 = 2000$, $x_2 = 5$, $price = \$500k$   one training example

size of the parameters $w_1, w_2$?

| $w_1 = 50$,  $w_2 = 0.1$,  $b = 50$ | $w_1 = 0.1$,  $w_2 = 50$,  $b = 50$ |
|---|---|
| | small   large |
| $\widehat{price} = \underbrace{50 * 2000}_{100,000K} + \underbrace{0.1 * 5}_{0.5k} + \underbrace{50}_{50K}$ | $\widehat{price} = \underbrace{0.1 * 2000k}_{200K} + \underbrace{50 * 5}_{150K} + \underbrace{50}_{50K}$ |
| $\widehat{price} = \$100,050.5\underline{k} = \$100,050,500$ | $\widehat{price} = \$500k$   more reasonable |

Stanford ONLINE   DeepLearning.AI   Andrew Ng

## Feature size and parameter size

|  | size of feature $x_j$ | size of parameter $w_j$ |
|---|---|---|
| size in feet² | ⟷ | ↔ |
| #bedrooms | ⟷ | ⟷ |

Features — scatterplot
$x_2$ # bedrooms
5, 0
$x_1$ size in feet²  300 ... 2000

Parameters — $J(\vec{w}, b)$ contour plot
$w_2$ # bedrooms
100, 10
$w_1$ size in feet²  0 ... 1

Stanford ONLINE   DeepLearning.AI   Andrew Ng

## Feature size and gradient descent

Features — scatterplot
$x_2$ # bedrooms
5, 0
$x_1$ size in feet²  300 ... 2000

$x_2$ # bedrooms rescaled
1, 0
$x_1$ size in feet² rescaled  0 ... 1

Parameters — contour plot $J(\vec{w}, b)$
$w_2$ # bedrooms
100, 10
$w_1$ size in feet²  0 ... 1

$w_2$ # bedrooms rescaled
$J(\vec{w}, b)$
$w_1$ size in feet² rescaled

Stanford ONLINE   DeepLearning.AI   Andrew Ng

## Feature scaling

$x_2$ # bedrooms

5
0

$300$ $x_1$ size in feet$^2$ $2000$

$x_2$ # bedrooms rescaled

1

0

$x_1$ size in feet$^2$ rescaled 1

$300 \leq x_1 \leq 2000$  $0 \leq x_2 \leq 5$

$$x_{1,scaled} = \frac{x_1}{2000} \quad \text{max}$$

$$x_{2,scaled} = \frac{x_2}{5} \quad \text{max}$$

$0.15 \leq x_{1,scaled} \leq 1$  $0 \leq x_{2,scaled} \leq 1$

---

## Mean normalization

$\mu_2 = 2.3$

average
$\mu_1 = 600$

$x_2$ # bedrooms

5
0

$300$ $x_1$ size in feet$^2$ $2000$

$x_2$ # bedrooms normalized

1

-1

1 $x_1$ size in feet$^2$ normalized

-1

$300 \leq x_1 \leq 2000$  $0 \leq x_2 \leq 5$

$$x_1 = \frac{x_1 - \mu_1}{2000-300} \quad \text{max-min}$$

$$x_2 = \frac{x_2 - \mu_2}{5-0} \quad \text{max-min}$$

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

---

## Z-score normalization

standard deviation $\sigma$

$\mu_2 = 2.3$

$x_2$ # bedrooms

$\sigma_1 = 450$
$\sigma_2 = 1.4$

$X_1$
$\mu_1$
$\sigma_1$

5
0

$300$ $\mu_1$ $2000$
$600$

$x_1$ size in feet$^2$

3

$x_2$ # bedrooms normalized

-3

3 $x_1$ size in feet$^2$ normalized

-3

$300 \leq x_1 \leq 2000$  $0 \leq x_2 \leq 5$

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

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

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

---

## Feature scaling

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

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

$0 \leq x_1 \leq 3$  okay, no rescaling

$-2 \leq x_2 \leq 0.5$  okay, no rescaling

$-100 \leq x_3 \leq 100$  too large → rescale

$-0.001 \leq x_4 \leq 0.001$  too small → rescale

$98.6 \leq x_5 \leq 105$  too large → rescale

# convergence of gradient descent

## Make sure gradient descent is working correctly

objective: $\min_{\vec{w},b} J(\vec{w},b)$

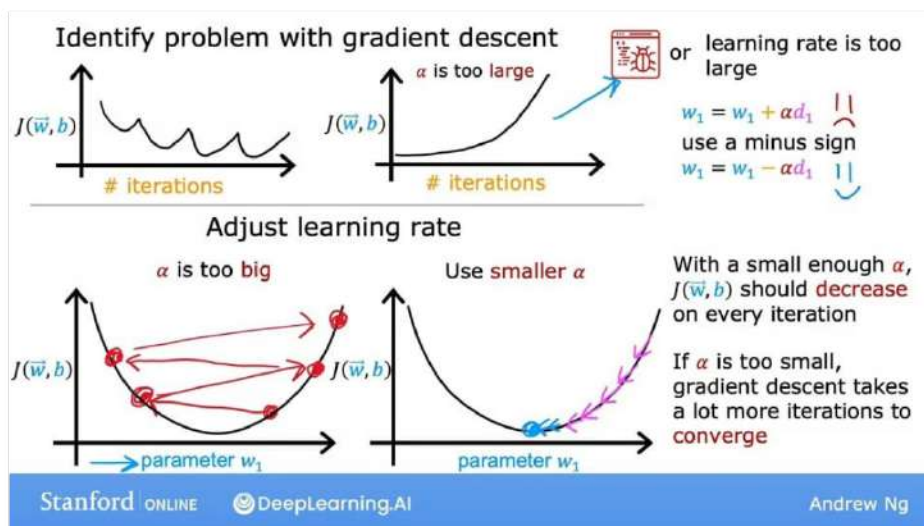$J(\vec{w},b)$ should **decrease** after every iteration

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

**learning curve**

$J(\vec{w},b)$ after **100** iterations

$J(\vec{w},b)$ after **200** iterations

$J(\vec{w},b)$ likely converged by **400** iterations

0   100   200   300   400

→ # iterations    $w, b$

\# iterations needed varies   30   1,000   100,000

**Automatic convergence test**
Let $\varepsilon$ "epsilon" be $10^{-3}$.
0.001

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

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

Stanford ONLINE    DeepLearning.AI    Andrew Ng

# choosing learning rate

## Identify problem with gradient descent

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

\# iterations

$\alpha$ is too large

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

\# iterations

or learning rate is too large

$w_1 = w_1 + \alpha d_1$
use a minus sign
$w_1 = w_1 - \alpha d_1$

## Adjust learning rate

$\alpha$ is too big

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

parameter $w_1$

Use smaller $\alpha$

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

parameter $w_1$

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

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

Stanford ONLINE    DeepLearning.AI    Andrew Ng

# Feature Engineering

## Feature engineering

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

frontage   depth

$area = frontage \times depth$

$x_3 = x_1 x_2$
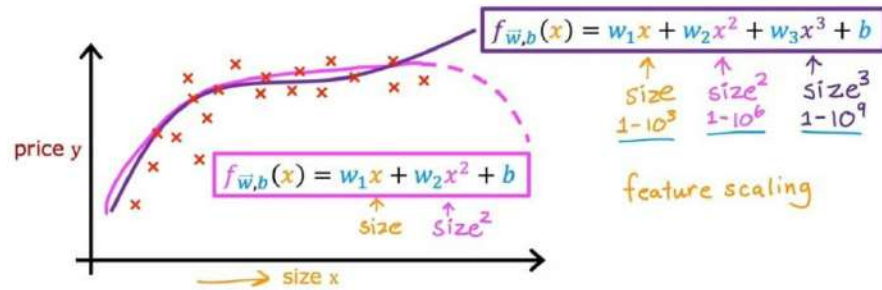new feature

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

frontage   depth

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

Stanford ONLINE    DeepLearning.AI    Andrew Ng

# Polynomial Regression

## Polynomial regression



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

size    size$^2$    size$^3$
$1-10^3$    $1-10^6$    $1-10^9$

feature scaling

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

size    size$^2$

price y

size x

## Choice of features



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

size    $\sqrt{size}$

price y

size x