

Machine Learning definition

- Arthur Samuel (1959). Machine Learning: Field of study that gives computers the ability to learn without being explicitly programmed.
- Tom Mitchell (1998) Well-posed Learning Problem: A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E.

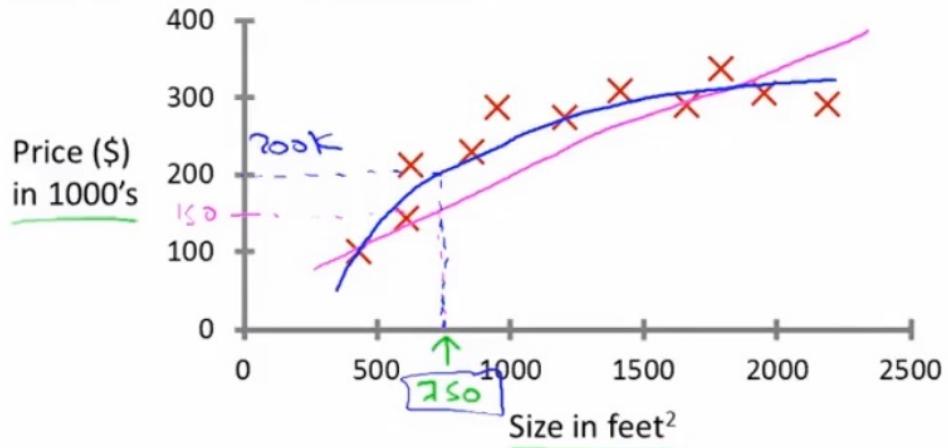
Task :

"A computer program is said to *learn* from experience E with respect to some task T and some performance measure P, if its performance on T, as measured by P, improves with experience E."

Suppose your email program watches which emails you do or do not mark as spam, and based on that learns how to better filter spam. What is the task T in this setting?

- Classifying emails as spam or not spam. T
- Watching you label emails as spam or not spam. E
- The number (or fraction) of emails correctly classified as spam/not spam. P
- None of the above—this is not a machine learning problem. P

Housing price prediction.



Supervised Learning

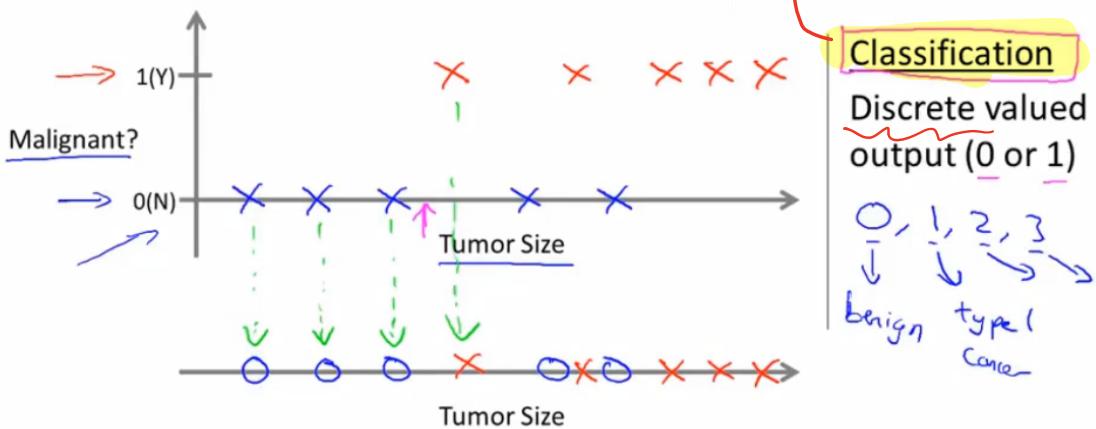
*right answers' given

Data Set

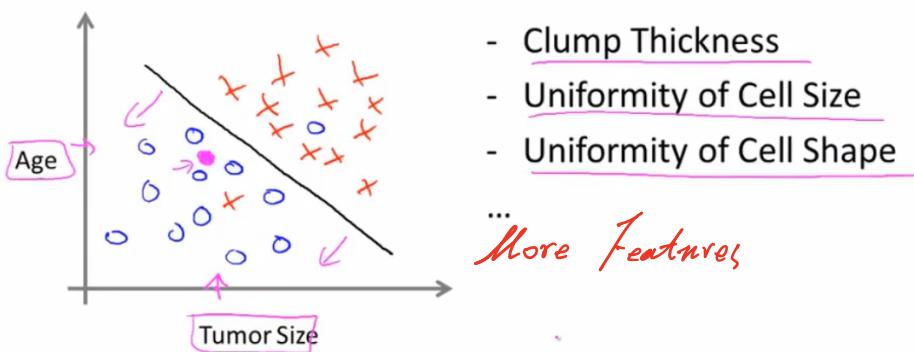
Regression: Predict continuous valued output (price)

Andrew Ng

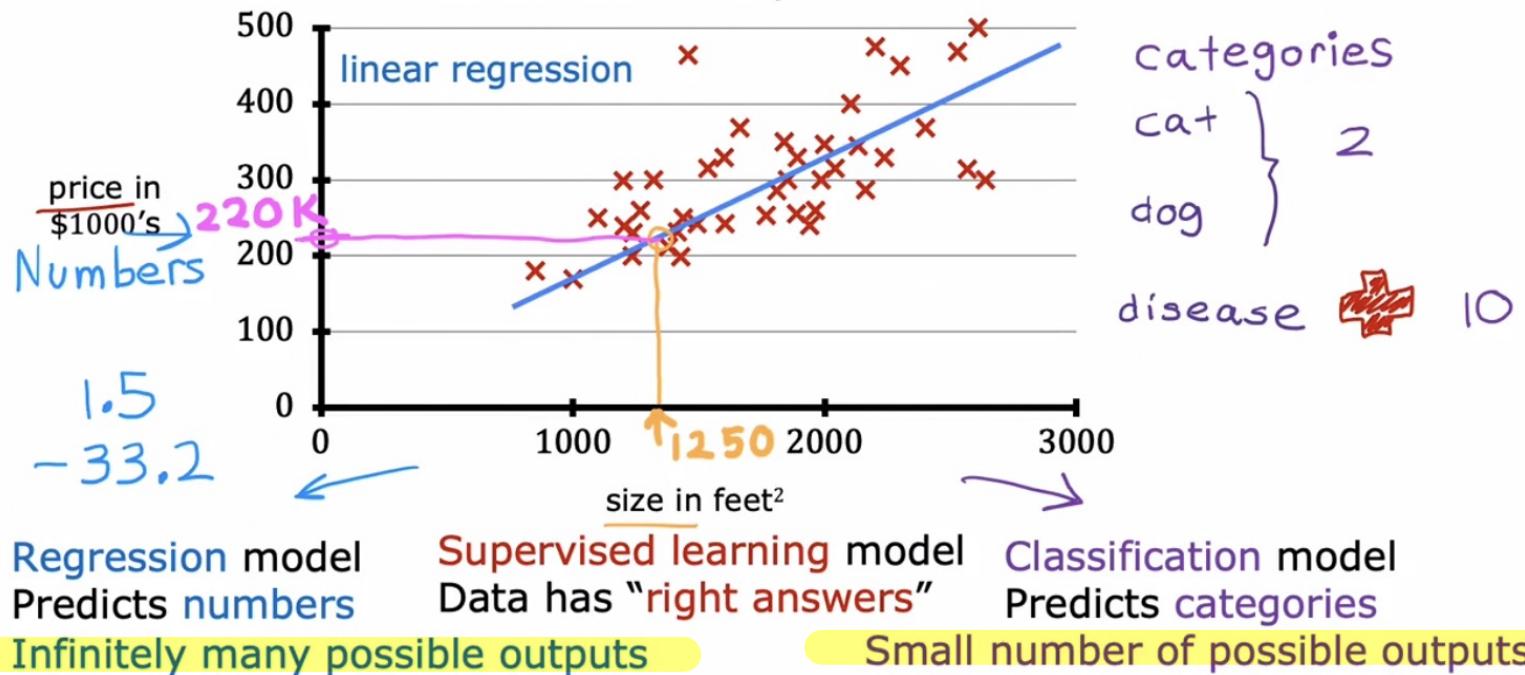
Breast cancer (malignant, benign)



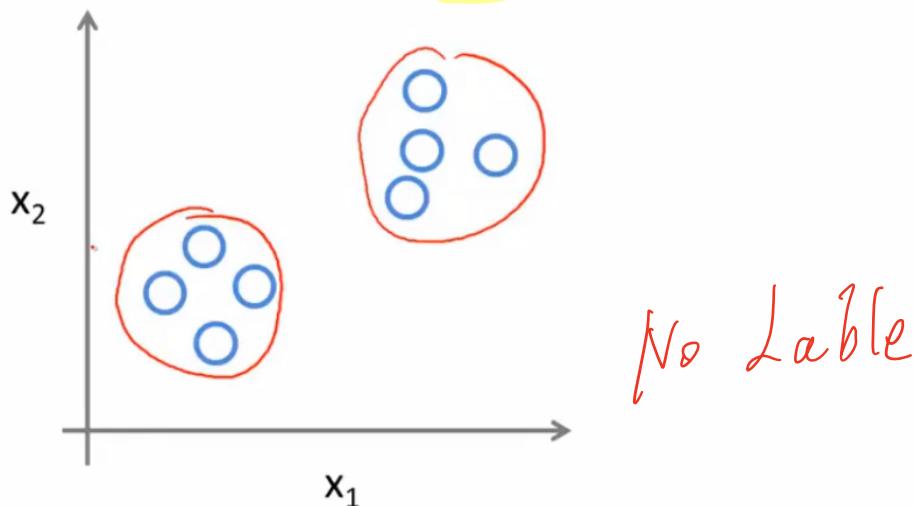
Andrew Ng



House sizes and prices



Unsupervised Learning



Andrew Ng

这就是聚类算法
So this is called a clustering algorithm.

Unsupervised learning

Data only comes with inputs x , but not output labels y .

Algorithm has to find structure in the data.

Clustering

Group similar data points together.

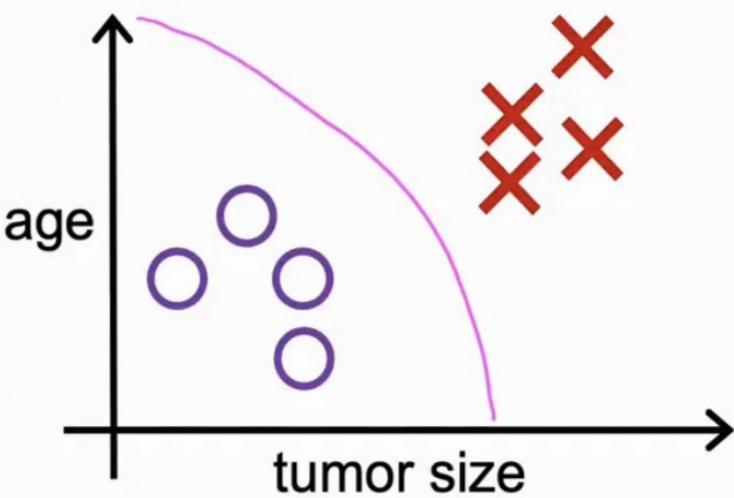
Dimensionality reduction

Compress data using fewer numbers.

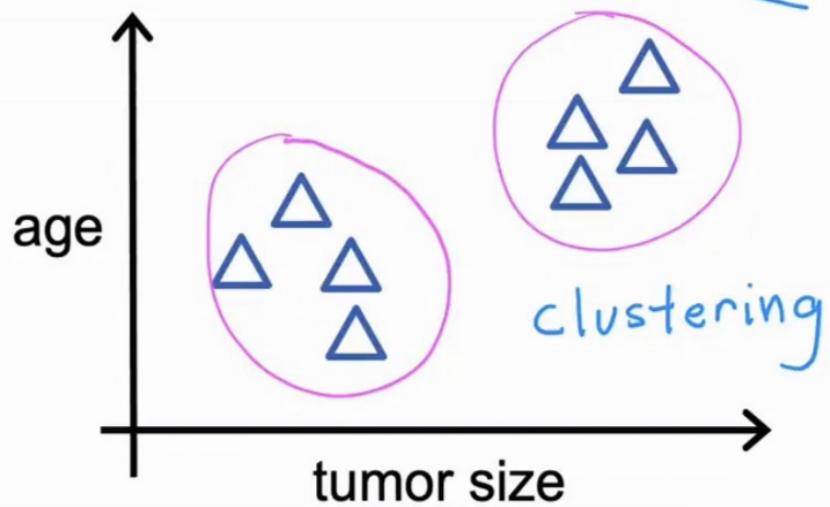
Anomaly detection

Find unusual data points.

Supervised learning
Learn from data **labeled**
with the “**right answers**”



Unsupervised learning
Find something interesting
in **unlabeled** data.



Training set of housing prices (Portland, OR)	Size in feet ² (x)	Price (\$) in 1000's(y)
→ 2104	460	
→ 1416	232	
→ 1534	315	
852	178	
...	...	

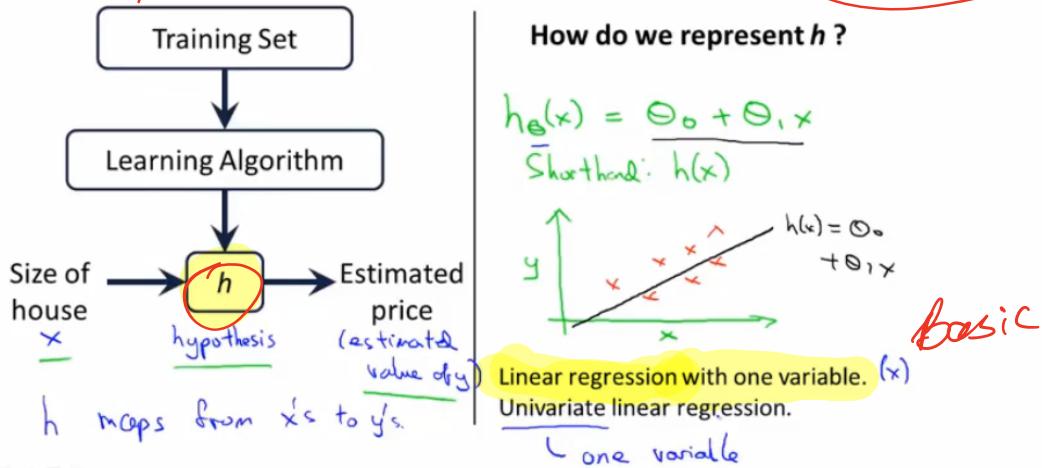
Notation:

- m = Number of training examples
- x 's = "input" variable / features
- y 's = "output" variable / "target" variable
- (x, y) - one training example
- $(x^{(i)}, y^{(i)})$ - i^{th} training example

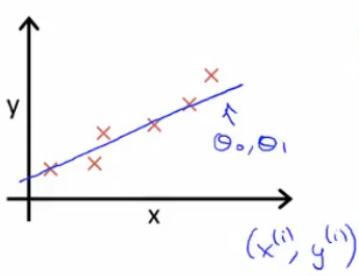
$$\begin{cases} x^{(1)} = 2104 \\ x^{(2)} = 1416 \\ \vdots \\ y^{(1)} = 460 \end{cases}$$

Andrew Ng

Supervised Learning Model



Cost Function



$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

\uparrow

$$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

Idea: Choose θ_0, θ_1 so that
 $h_\theta(x)$ is close to y for our
training examples (x, y)

$$\text{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

Cost function

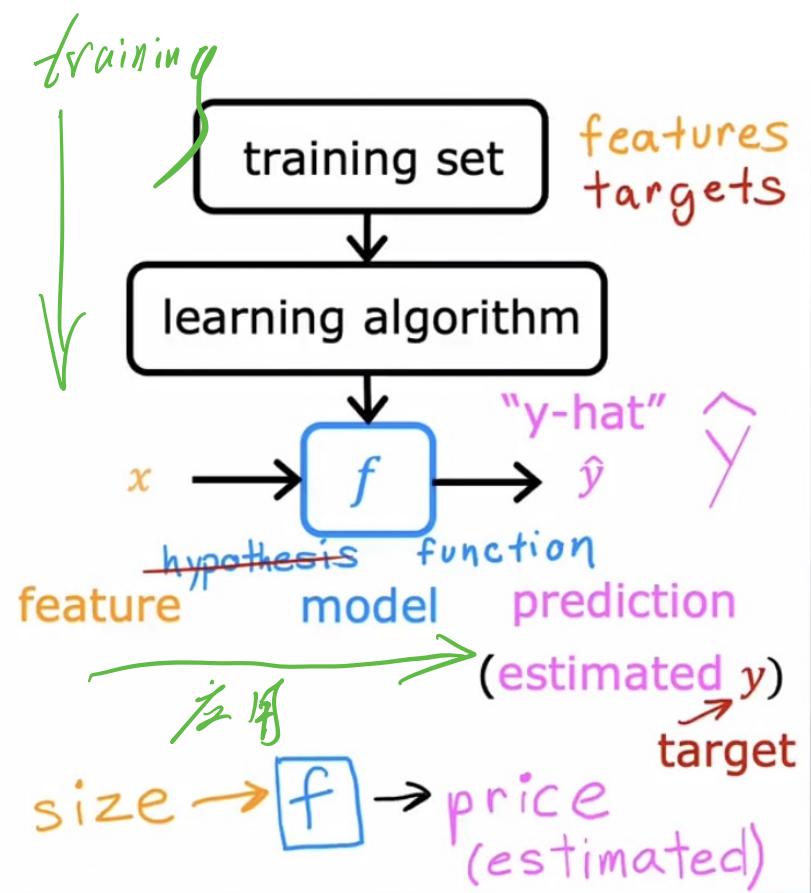
Squared error function

Andrew Ng

error cost function is probably the most
commonly used one for regression problems.

但是平方误差代价函数可能是解决回归问题最常用的手段了

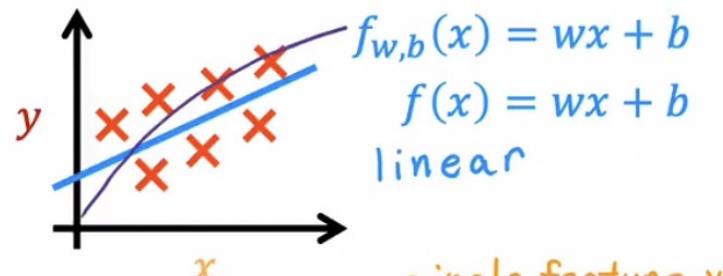
Regression



How to represent f ?

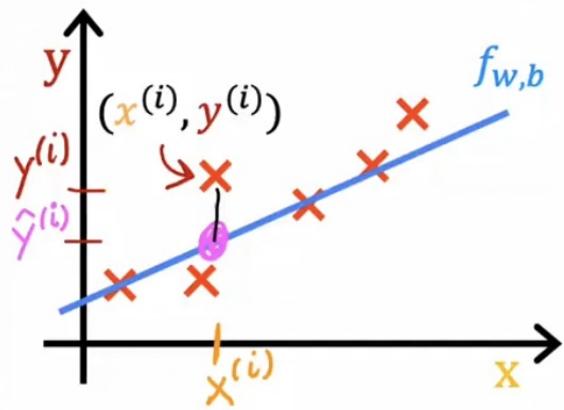
$$f_{w,b}(x) = wx + b$$

$$f(x)$$



single feature x
Linear regression with one variable.
size

Univariate linear regression.
one variable



$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

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

Cost function: Squared error cost function

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

m = number of training examples

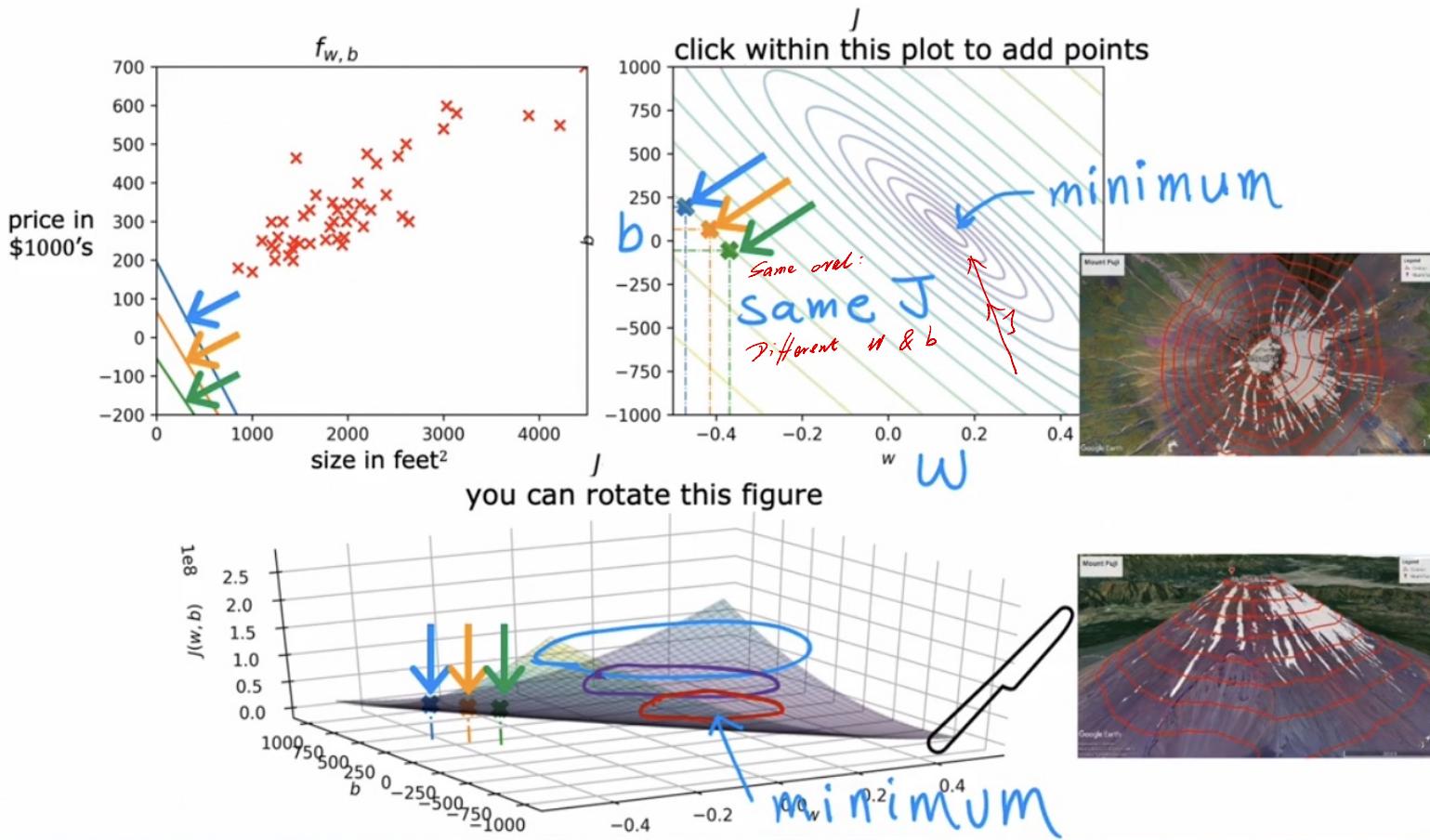
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Find w, b :

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

$\min J(w, b)$ 返回 J^*

$\arg \min J(w, b)$ 返回 w^*, b^*



To find $w \& b$ that $\min J$

Gradient Descent
梯度下降

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

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

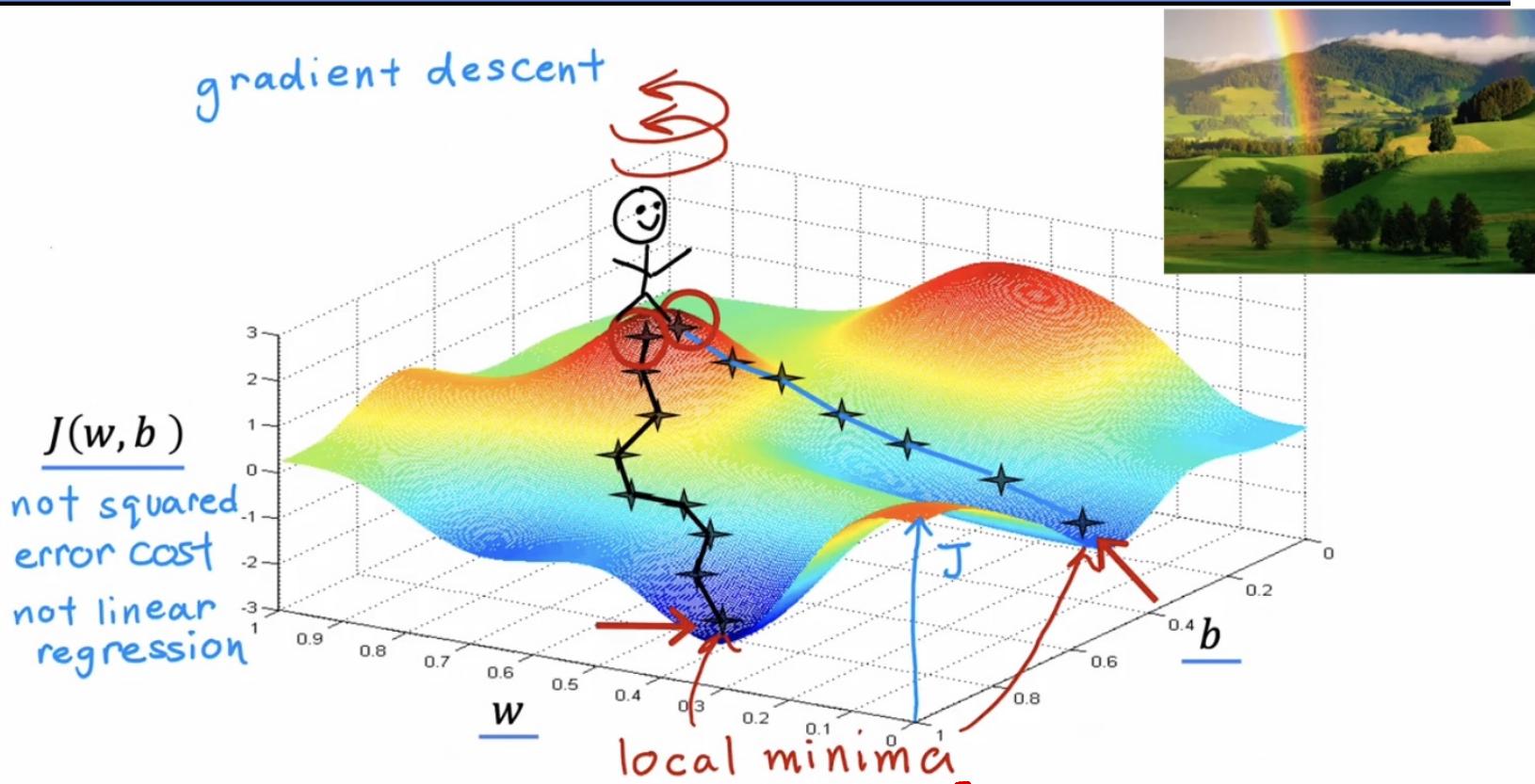
Outline:

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

Keep changing w, b to reduce $J(w, b)$

Until we settle at or near a minimum

may have >1 minimum



全局极小
Global minimum: 全局最小

Gradient descent algorithm

Repeat until convergence

$$\left\{ \begin{array}{l} w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ b = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{array} \right.$$

$\alpha \sim (0, 1)$
Learning rate
Derivative

Simultaneously update w and b

Assignment

$a = c$
 $a = a + 1$
Code

Truth assertion

$a = c$
 $a = a + 1$
Math
 $a == c$

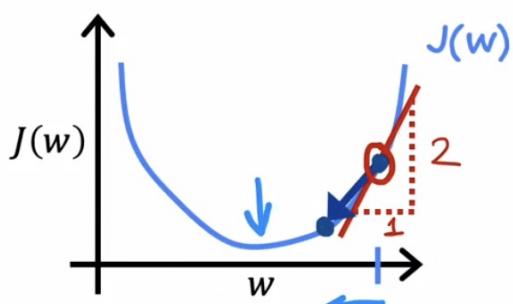
Correct: Simultaneous update

$$\left\{ \begin{array}{l} tmp_w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ tmp_b = b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w = tmp_w \\ b = tmp_b \end{array} \right. \quad \text{同步更新}$$

Incorrect

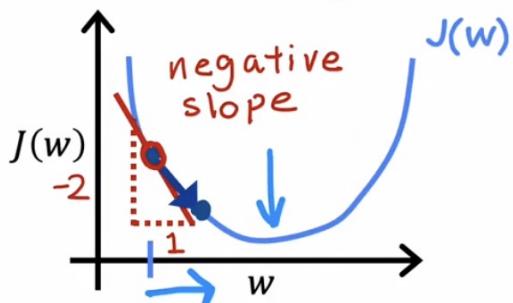
$$\left\{ \begin{array}{l} tmp_w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ w = tmp_w \\ tmp_b = b - \alpha \frac{\partial}{\partial b} J(w, b) \\ b = tmp_b \end{array} \right.$$

Derivative



if $b > 0$:

$$w = w - \alpha \frac{\partial}{\partial w} J(w) > 0$$



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

$$\frac{\partial}{\partial w} J(w) < 0$$

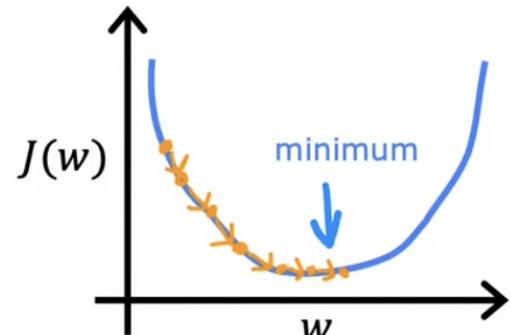
$w = w - \alpha \cdot (\text{negative number})$

The Learning Rate controls how big of a step you take when updating the model's parameters.

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

If α is too small...

Gradient descent may be slow.

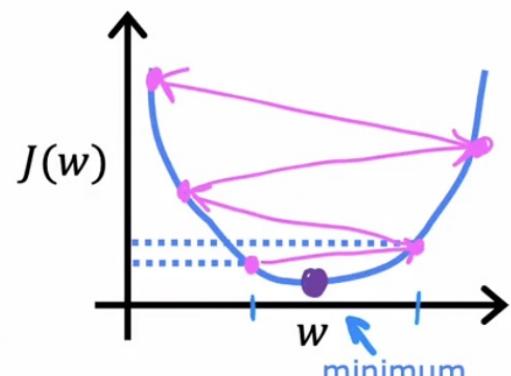


If α is too large...

Gradient descent may:

- Overshoot, never reach minimum
- Fail to converge, diverge

~~无法收敛，发散~~



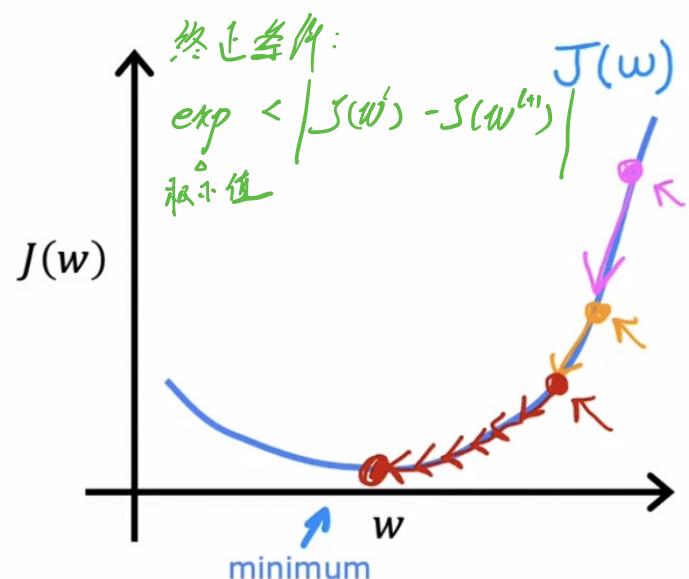
Can reach local minimum with fixed learning rate α

$w = w - \alpha \frac{d}{dw} J(w)$	smaller not as large large
------------------------------------	----------------------------------

Near a local minimum,

- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate α



Linear regression model

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

Cost function

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

Gradient descent algorithm

repeat until convergence {

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

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

}

未停止

“Batch” gradient descent

“Batch”: Each step of gradient descent uses all the training examples.

other gradient
descent: subsets

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

$m = 47$

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

多元(特征)线性回归

Multiple features (variables)

	Size in feet ² x_1	Number of bedrooms x_2	Number of floors x_3	Age of home in years x_4	Price (\$) in \$1000's
$i=2$	2104	5	1	45	460
	1416	3	2	40	232
	1534	3	2	30	315
	852	2	1	36	178

$x_j = j^{th}$ 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

$j=1 \dots 4$

$n=4$

$$\vec{x}^{(2)} = [1416 \ 3 \ 2 \ 40]$$

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

特征向量

有时为了强调 x_2 不是一个数字 它实际上是一个向量的数字列表

actually a list of numbers that is a vector,

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

$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n]$ parameters
of the model

b is a number

vector $\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$

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

dot product multiple linear regression
 ↗
 (not multivariate regression)

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

```
w = np.array([1, 0, 2, 5, -3, 3])
```

```
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

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

`range(0,n) → j = 0...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
```

Dot product



向量化的概念 & 矩阵

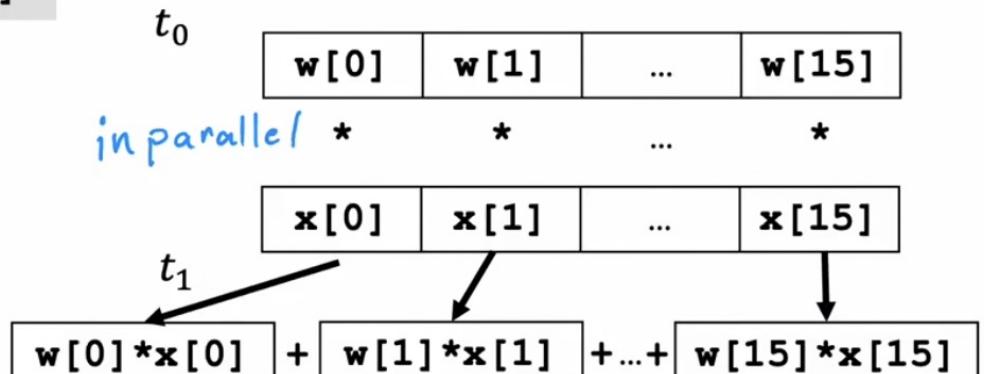
Without vectorization

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

$t_0 \quad f + w[0] * x[0]$
 $t_1 \quad f + w[1] * x[1]$
 \dots
 $t_{15} \quad f + w[15] * x[15]$

Vectorization

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



efficient \rightarrow scale to large datasets

Gradient descent

$\vec{w} = (w_1 \quad w_2 \quad \dots \quad w_{16})$ ~~parameters~~
 derivatives $\vec{d} = (d_1 \quad d_2 \quad \dots \quad d_{16})$

```
w = np.array([0.5, 1.3, ... 3.4])
d = np.array([0.3, 0.2, ... 0.4])
```

compute $w_j = w_j - \text{learning rate } \alpha d_j$ for $j = 1 \dots 16$

Without vectorization

$$\begin{aligned} w_1 &= w_1 - 0.1d_1 \\ w_2 &= w_2 - 0.1d_2 \\ &\vdots \\ w_{16} &= w_{16} - 0.1d_{16} \end{aligned}$$

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

With vectorization

The diagram illustrates the computation of the gradient descent update $\vec{w} = \vec{w} - 0.1\vec{d}$. It shows two vectors, w and d , represented as horizontal arrays of 16 elements each. The vector w is being updated by subtracting $0.1 \cdot d$. The diagram shows the subtraction operation as the difference between the two vectors, with arrows indicating the element-wise subtraction and the scalar multiplication by 0.1.

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

Gradient Descent for Multiple Regression

Previous notation

Parameters w_1, \dots, w_n
 b

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

Cost function $J(\underbrace{w_1, \dots, w_n}_b)$

Vector notation

\vec{w} ↗ vector of length n
 $\vec{w} = [w_1 \dots w_n]$
 b still a number
 $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
 $J(\vec{w}, b)$ ↗ dot product

Gradient descent

```
repeat {
     $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\underbrace{w_1, \dots, w_n}_b, b)$ 
     $b = b - \alpha \frac{\partial}{\partial b} J(\underbrace{w_1, \dots, w_n}_b, b)$ 
}
```

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

Gradient descent

One feature

```
repeat {
     $\vec{w} = \vec{w} - \alpha \left( \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \vec{x}^{(i)} \right)$ 
    ↳  $\frac{\partial}{\partial \vec{w}} J(\vec{w}, b)$ 
     $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$ 
    simultaneously update  $\vec{w}, b$ 
}
```

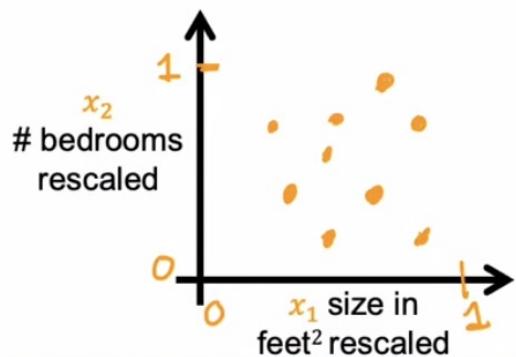
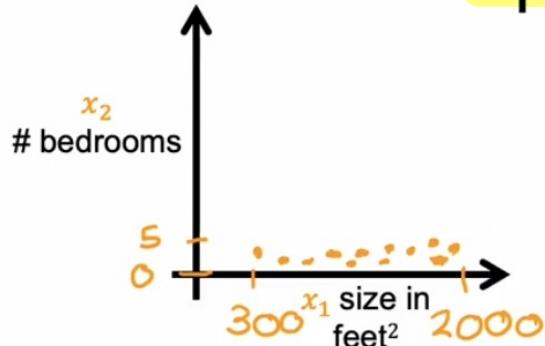
n features ($n \geq 2$)

```
repeat {
     $j=1$   $w_1 = w_1 - \alpha \left( \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \vec{x}_1^{(i)} \right)$ 
    :
     $j=n$   $w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \vec{x}_n^{(i)}$ 
     $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$ 
    simultaneously update
     $w_j$  (for  $j = 1, \dots, n$ ) and  $b$ 
}
```

Make Gradient Descent Faster

Feature scaling

缩放

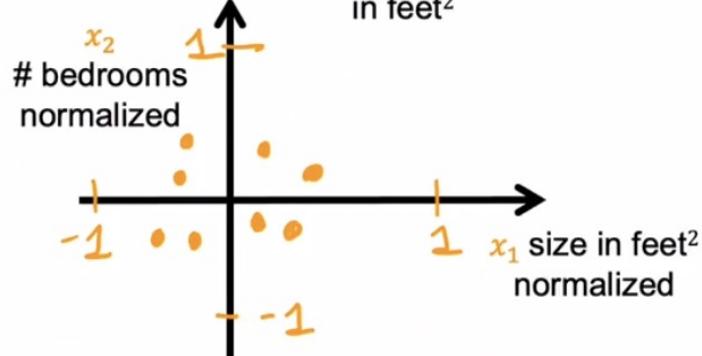
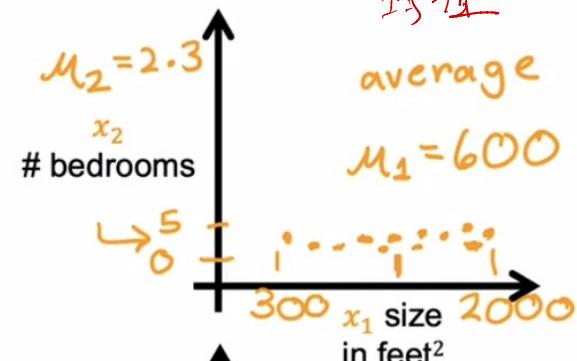


$$300 \leq x_1 \leq 2000$$
$$x_{1,scaled} = \frac{x_1}{2000}$$
$$0.15 \leq x_{1,scaled} \leq 1$$

$$0 \leq x_2 \leq 5$$
$$x_{2,scaled} = \frac{x_2}{5}$$
$$0 \leq x_{2,scaled} \leq 1$$

Mean normalization

均值



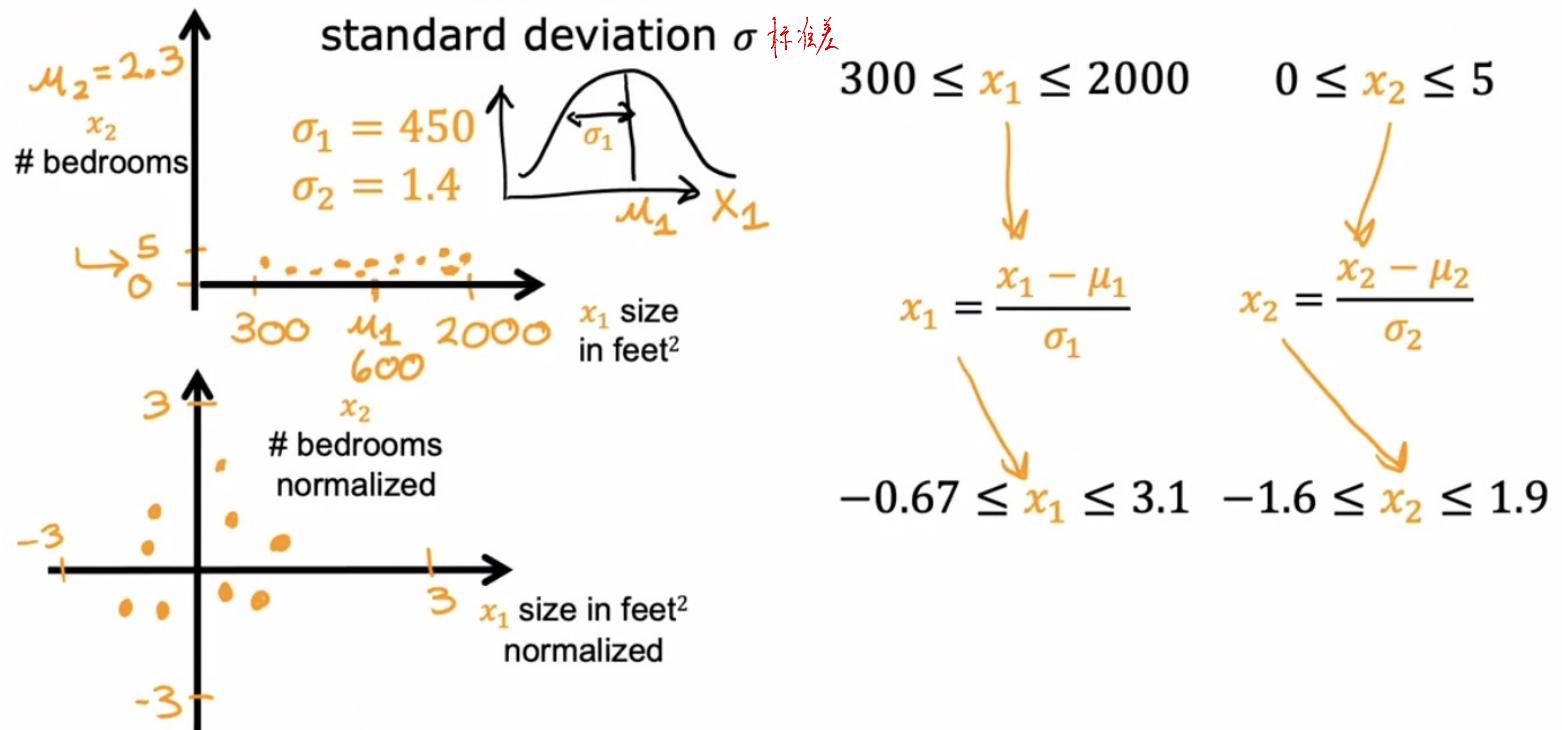
$$300 \leq x_1 \leq 2000$$

$$x_1 = \frac{x_1 - \mu_1}{2000 - 300}$$
$$-0.18 \leq x_1 \leq 0.82$$

$$0 \leq x_2 \leq 5$$

$$x_2 = \frac{x_2 - \mu_2}{5 - 0}$$
$$-0.46 \leq x_2 \leq 0.54$$

Z-score normalization



Feature scaling

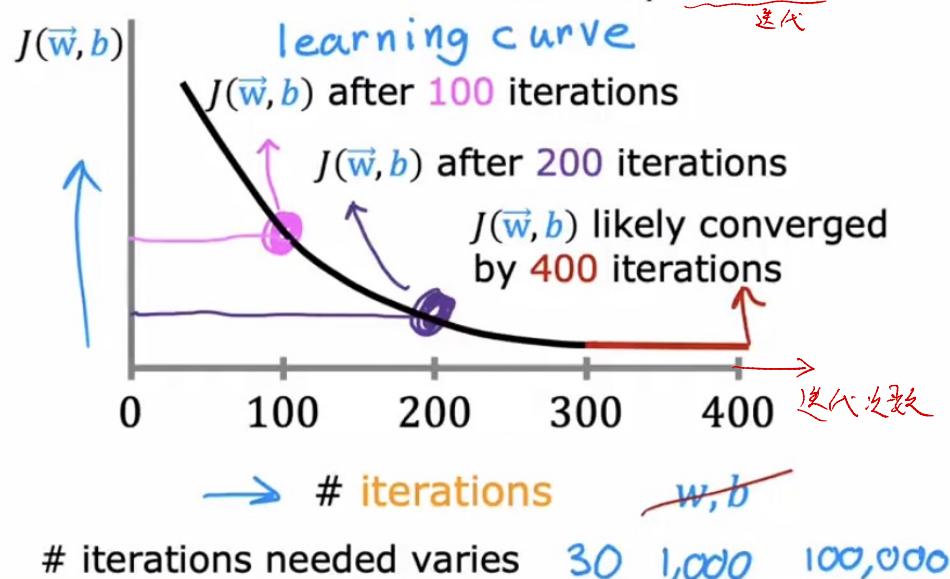
aim for about $-1 \leq x_j \leq 1$ for each feature x_j

$-3 \leq x_j \leq 3$	}	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 \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

Make sure gradient descent is working correctly

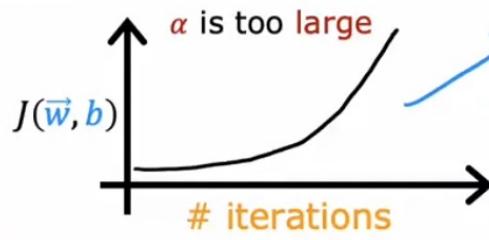
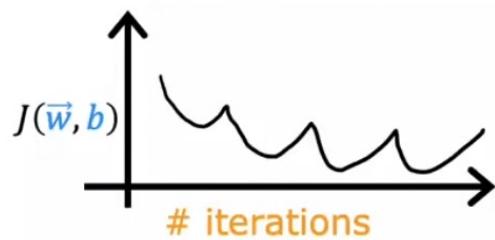
objective: $\min_{\vec{w}, b} J(\vec{w}, b)$ $J(\vec{w}, b)$ should decrease after every iteration



Automatic convergence test
Let ε "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)

Identify problem with gradient descent

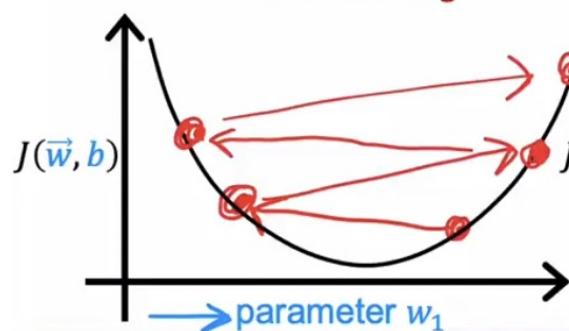


or learning rate is too large

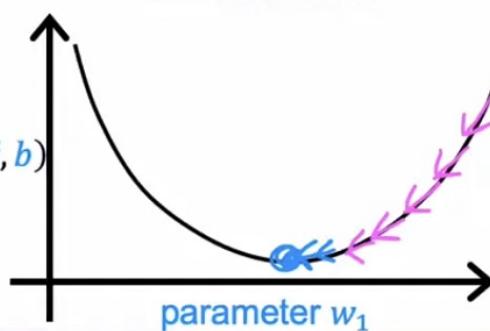
$$\begin{aligned} w_1 &= w_1 + \alpha d_1 \\ &\text{use a minus sign} \\ w_1 &= w_1 - \alpha d_1 \end{aligned}$$

Adjust learning rate

α is too big



Use smaller α



With a small enough α , $J(\vec{w}, b)$ should decrease on every iteration
If not, code error/bug

If α is too small, gradient descent takes a lot more iterations to converge

Feature engineering

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

frontage depth

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

$$x_3 = x_1 x_2$$

new feature

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



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

Classification

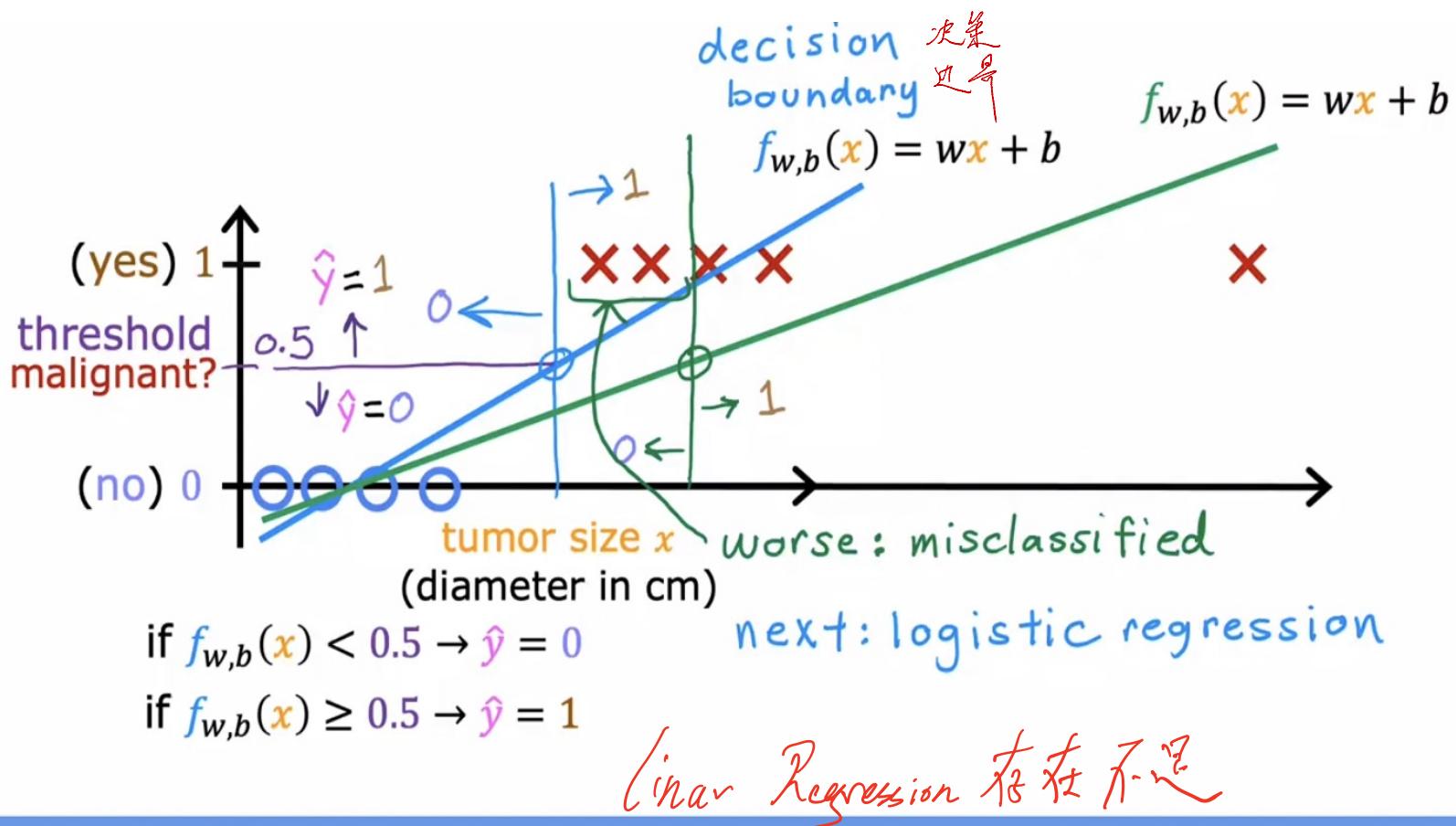
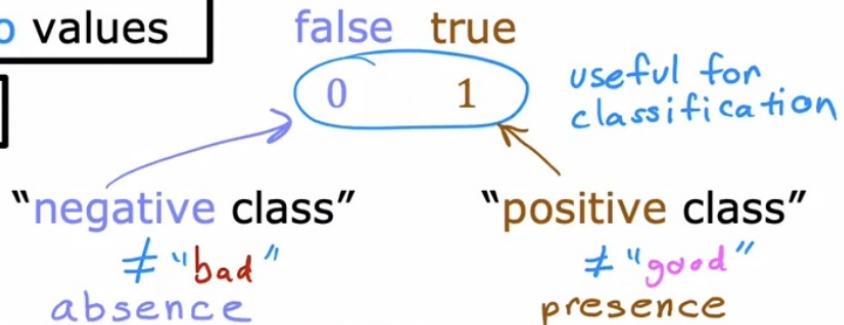
Classification

Question	Answer " y "	
Is this email <u>spam</u> ?	no	yes
Is the transaction <u>fraudulent</u> ?	no	yes
Is the tumor <u>malignant</u> ?	no	yes

y can only be one of two values

"binary classification"

class = category

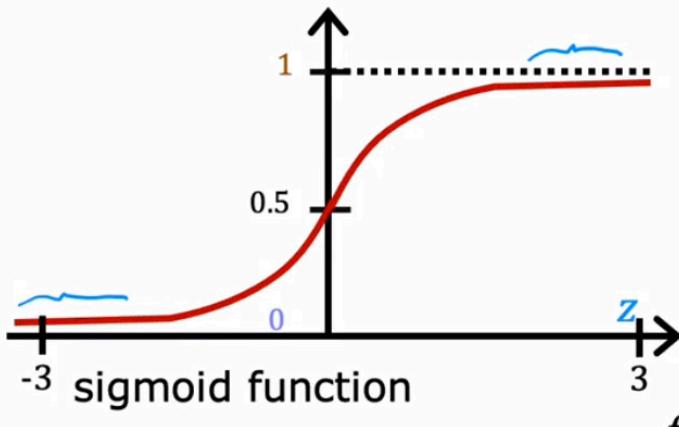


顺便说一下，逻辑回归这个名字让人困惑的一点是，尽管它里面有回归 (regression) 这个词，但它实际上是用来分类的。

By the way one thing confusing about the name logistic regression is that even though it has the word of regression in it is actually used for classification.

Logistic Regression

Want outputs between 0 and 1



outputs between 0 and 1

$$g(z) = \frac{1}{1+e^{-z}} \quad 0 < g(z) < 1$$

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

$$z = \vec{w} \cdot \vec{x} + b$$

\downarrow
 \downarrow

$$g(z) = \frac{1}{1+e^{-z}}$$

$$f_{\vec{w}, b}(\vec{x}) = g(\underbrace{\vec{w} \cdot \vec{x} + b}_z) = \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}}$$

"logistic regression" $e \approx 2.7$

Interpretation of logistic regression output

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}}$$

"probability" that class is 1

$$f_{\vec{w}, b}(\vec{x}) = P(y = 1 | \vec{x}; \vec{w}, b)$$

Probability that y is 1,
given input \vec{x} , parameters \vec{w}, b

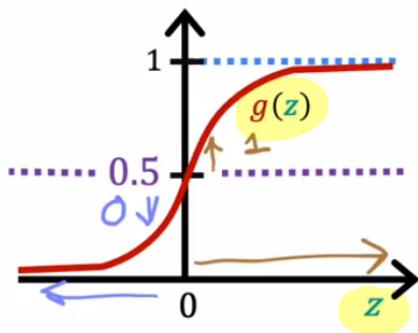
Example:

- x is "tumor size"
- y is 0 (not malignant)
or 1 (malignant)

$$P(y = 0) + P(y = 1) = 1$$

$$f_{\vec{w}, b}(\vec{x}) = 0.7$$

70% chance that y is 1



$$f_{\vec{w}, b}(\vec{x}) = g(\vec{w} \cdot \vec{x} + b) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

z 是中间值, 可自由设置
 \hat{y} 是预测值

0 or 1? threshold

Is $f_{\vec{w}, b}(\vec{x}) \geq 0.5?$

Yes: $\hat{y} = 1$

No: $\hat{y} = 0$

When is $f_{\vec{w}, b}(\vec{x}) \geq 0.5?$

Sigmoid

$g(z) \geq 0.5$

z 、 $g(z)$ 均可作
判定条件.

$z \geq 0$

$\vec{w} \cdot \vec{x} + b \geq 0$

$\hat{y} = 1$

$\vec{w} \cdot \vec{x} + b < 0$

$\hat{y} = 0$

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

$$z = \vec{w} \cdot \vec{x} + b$$

\downarrow

\downarrow

$$g(z) = \frac{1}{1+e^{-z}}$$

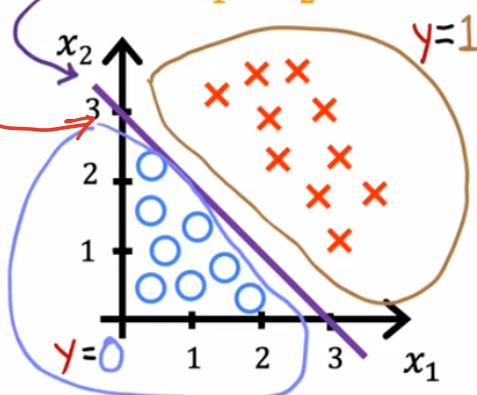
Decision boundary

$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + b) = g(1 \cdot x_1 + 1 \cdot x_2 - 3)$$

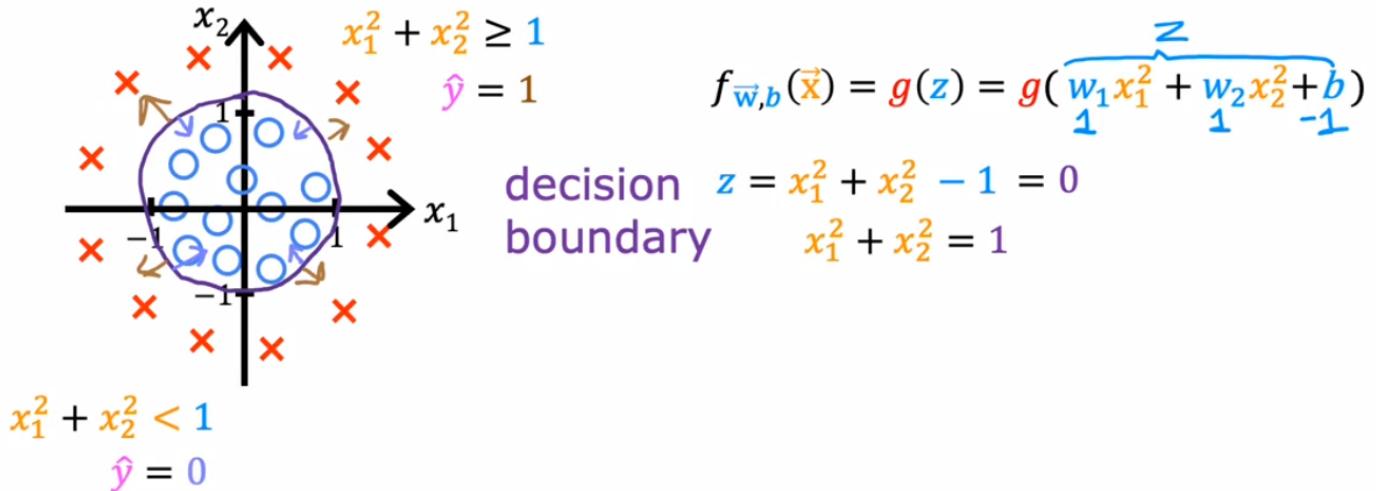
Decision boundary $z = \vec{w} \cdot \vec{x} + b = 0 \Rightarrow g(z) = 0.5$ 多维特征.

$$z = x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3$$



Non-linear decision boundaries



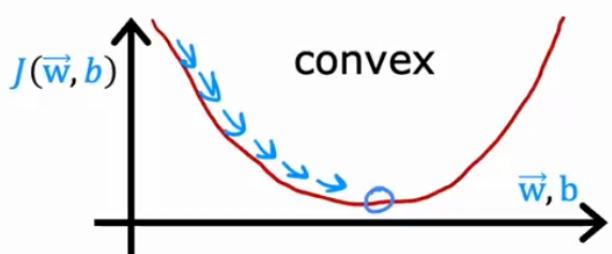
Cost Function

Squared error cost *Bad Choice*

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

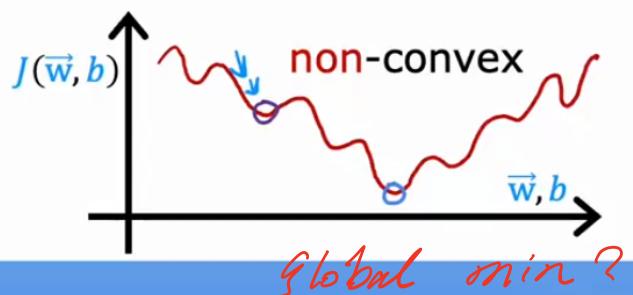
linear regression

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



logistic regression

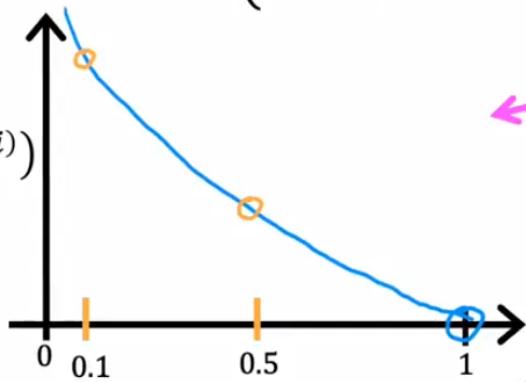
$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$



Logistic loss function

Part of Cost Function

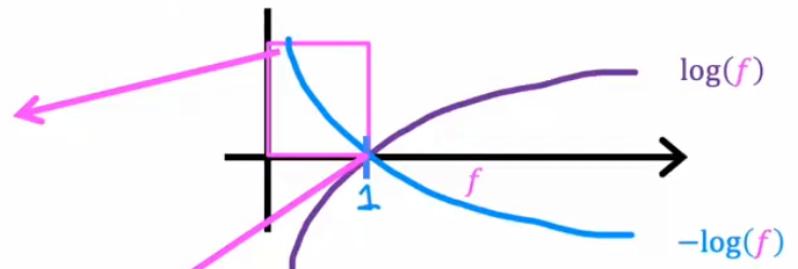
$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$



if $y^{(i)} = 1$

As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 1$ then loss $\rightarrow 0$

As $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 0$ then loss $\rightarrow \infty$



Loss is lowest when
 $f_{\vec{w}, b}(\vec{x}^{(i)})$ predicts
close to true label $y^{(i)}$.

Simplified cost function

$$\text{loss} \quad L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

$$\text{cost} \quad J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})]$$

Convex
(single global minimum)

$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))]$$

maximum likelihood
(don't worry about it!)

Gradient descent

cost

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

repeat {

$j = 1 \dots n$

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

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

} simultaneous updates

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

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

Gradient descent for logistic regression

repeat {

looks like linear regression!

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

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

} simultaneous updates

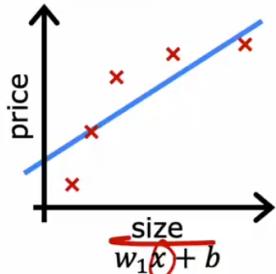
- Same concepts:
- Monitor gradient descent (learning curve)
 - Vectorized implementation
 - Feature scaling

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

Logistic regression $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{(-\vec{w} \cdot \vec{x} + b)}}$

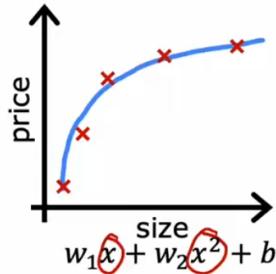
Overfitting & Underfitting

Regression example



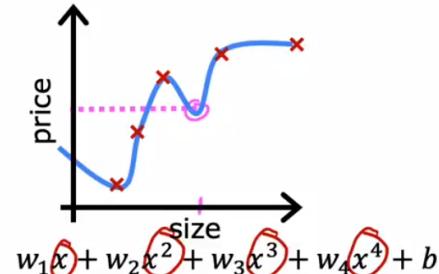
- Does not fit the training set well

high bias



- Fits training set pretty well

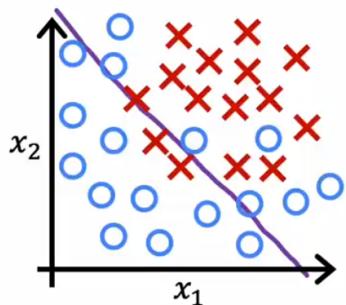
generalization
沒學到



- Fits the training set extremely well

high variance

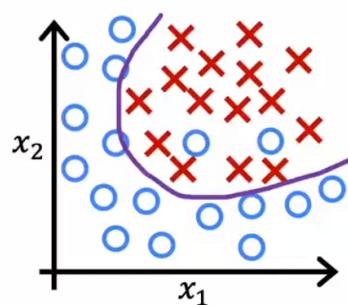
Classification



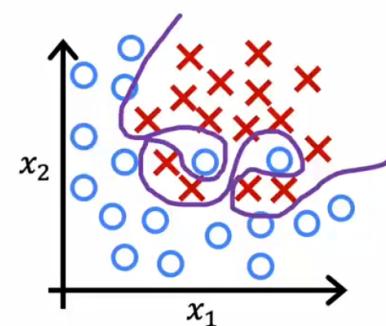
$$z = w_1 x_1 + w_2 x_2 + b$$

$$f_{\vec{w}, b}(\vec{x}) = g(z)$$

g is the sigmoid function
high bias
高偏差



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_2^2 + w_5 x_1 x_2 + b$$



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 + w_5 x_1^2 x_3 + w_6 x_1^3 x_2 + \dots + b$$

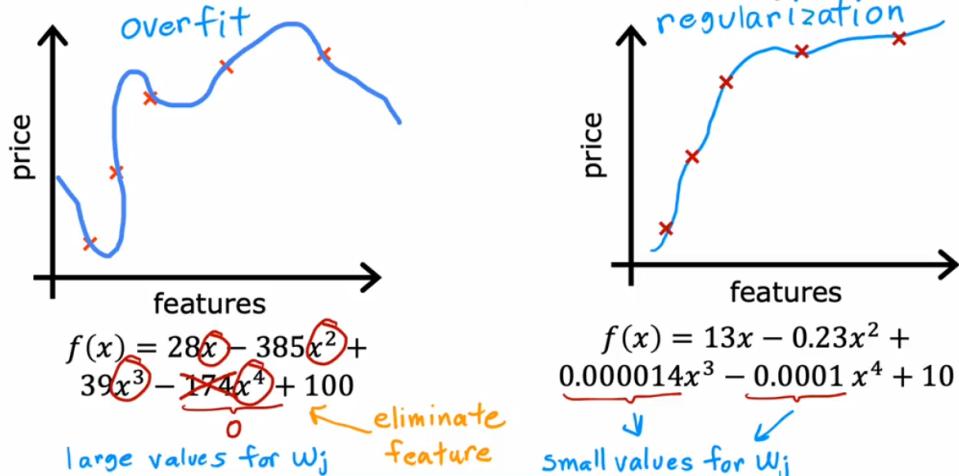
Addressing overfitting

Options

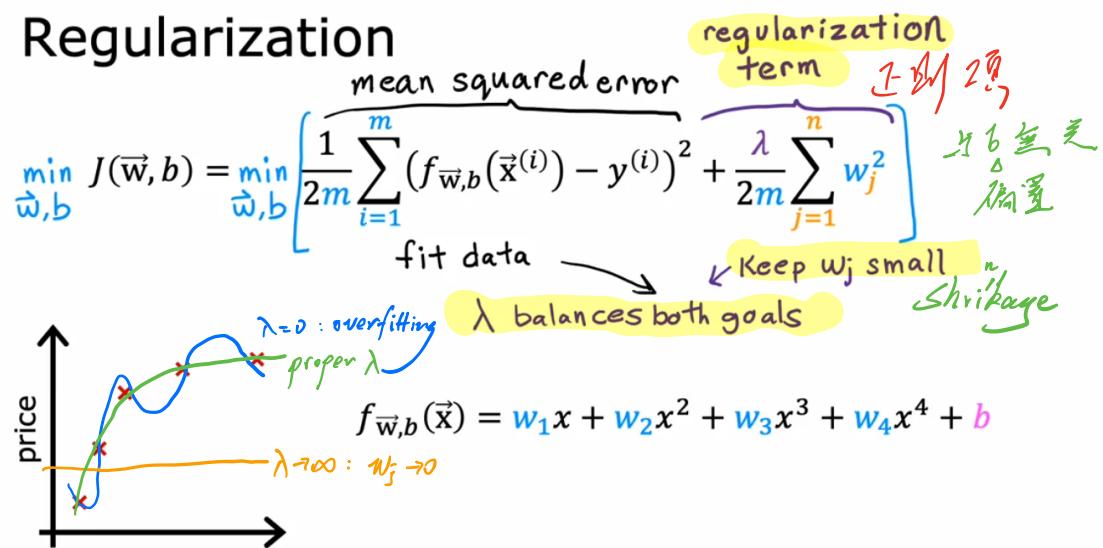
1. Collect more data
2. Select features
 - Feature selection *in course 2*
3. Reduce size of parameters
 - "Regularization" *next videos!*
± w_j f(x)

Regularization

Reduce the size of parameters w_j



Regularization



Regularized linear regression

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right]$$

Gradient descent

repeat {

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

$j = 1, \dots, n$

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

} simultaneous update

$= \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j$

$= \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$ don't have to regularize b

Implementing gradient descent

repeat {

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

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

} simultaneous update $j = 1 \dots n$

$$w_j = \underbrace{w_j - \alpha \frac{\lambda}{m} w_j}_{w_j \left(1 - \alpha \frac{\lambda}{m}\right)} - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{usual update}}$$

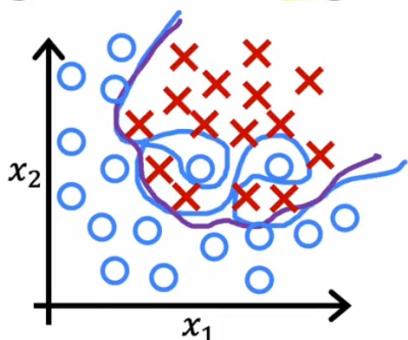
$$\alpha \frac{\lambda}{m}$$

$$0.01 \frac{1}{50} = 0.0002$$

$$w_j \left(1 - 0.0002\right)$$

$$0.9998$$

Regularized logistic regression



$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 + w_5 x_1^2 x_2^3 + \dots + b$$

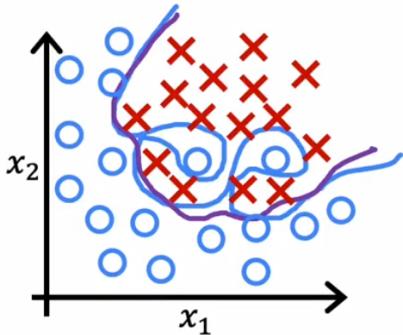
$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-z}}$$

Cost function

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$\min_{\vec{w}, b} J(\vec{w}, b) \rightarrow w_j \downarrow$

Regularized logistic regression



$$z = w_1x_1 + w_2x_2 + w_3x_1^2x_2 + w_4x_1^2x_2^2 + w_5x_1^2x_2^3 + \dots + b$$
$$f_{\vec{w},b}(\vec{x}) = \frac{1}{1 + e^{-z}}$$

Cost function

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$\min_{\vec{w}, b} J(\vec{w}, b) \rightarrow w_j \downarrow$

有一点需要注意，“线性回归”中的“线性”指的是因变量关于参数是线性的，而不是关于自变量，比如模型： $Y_i = \beta_1 X_i + \beta_2 X_i^2 + \varepsilon_i$ ，因变量对于参数 β_1 和 β_2 是线性的，但对于自变量 X_i 来说就是非线性的，这仍然是一个线性回归模型。

不过对于线性回归，它们的代价函数图像并非上面那样，线性回归的代价函数图像总是**一个凸函数 (convex function)**，如下图，它并没有“局部最优解”，它只有一个“全局最优解”。

