





The dataset that is used to train the model is called a training set.

**Notation:**

$x$  = “input” variable  
feature

$y$  = “output” variable  
“target” variable

$m$  = number of training examples

( $x$ ,  $y$ ) = single training example

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

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

(1<sup>st</sup>, 2<sup>nd</sup>, 3<sup>rd</sup> ...)

Cost function formula:

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

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

$$f_{w,b}(x^{(i)}) = wx^{(i)} + b \quad \text{Find } w, b: \\ \hat{y}^{(i)} \text{ is close to } y^{(i)} \text{ for all } (x^{(i)}, y^{(i)}).$$

Cost function: Squared error cost function

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

$m$  = 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



## Gradient descent algorithm

Repeat until convergence

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

Learning rate  
Derivative

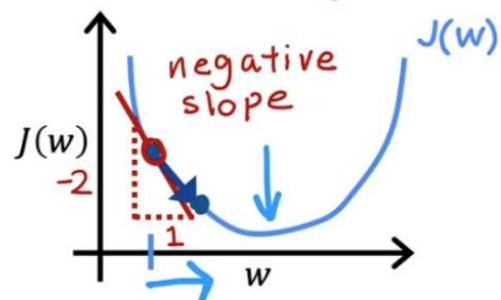
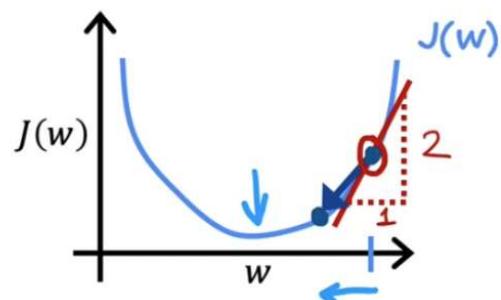
Simultaneously  
update w and b

**Correct: Simultaneous update**

$$\begin{aligned} \text{tmp\_w} &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \text{tmp\_b} &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w &= \text{tmp\_w} \\ b &= \text{tmp\_b} \end{aligned} \quad \left. \right\}$$

**Incorrect**

$$\begin{aligned} \text{tmp\_w} &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ w &= \text{tmp\_w} \\ \text{tmp\_b} &= b - \alpha \frac{\partial}{\partial b} J(\text{tmp\_w}, b) \\ b &= \text{tmp\_b} \end{aligned}$$



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

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

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

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

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

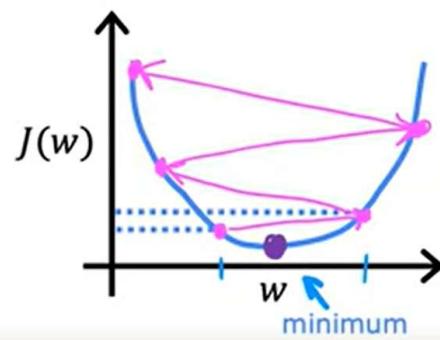
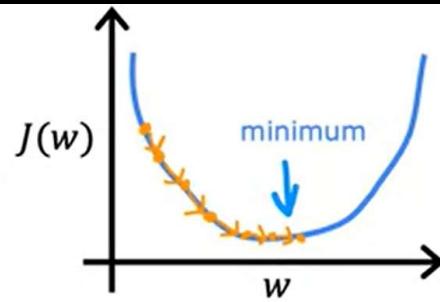
If  $\alpha$  is too small...

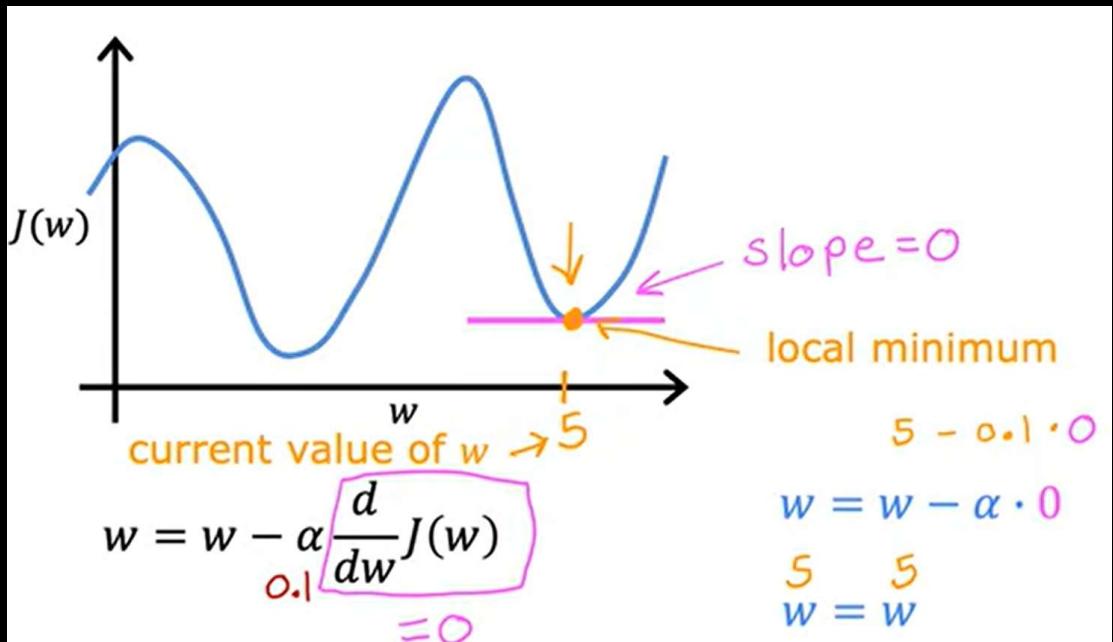
Gradient descent may be slow.

If  $\alpha$  is too large...

Gradient descent may:

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





Near a local minimum,

- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate  $\alpha$

(Optional)

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

$$= \cancel{\frac{1}{2m}} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)}) \cancel{2X^{(i)}} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}}$$

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

$$= \cancel{\frac{1}{2m}} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)}) \cancel{2} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})}$$

no  $x^{(i)}$

## Gradient descent algorithm

```
repeat until convergence {  
     $w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$   
     $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$   
}
```

	$x_1$	$x_2$	$x_3$	$x_4$	
$i=2$	2104	5	1	45	
	1416	3	2	40	
	1534	3	2	30	
	852	2	1	36	
	...	...	...	...	

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

$x_j^{(i)}$  = value of feature  $j$  in the  $i^{\text{th}}$  training example

$\vec{x}^{(i)}$  = the column vector of all the feature inputs of the  $i^{\text{th}}$  training example

$m$  = the number of training examples

$n = |\vec{x}^{(i)}|$ ; (the number of features)

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

$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n]$  parameters  
 $b$  is a number  
 vector  $\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$

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

$\uparrow$  dot product      multiple linear regression

Parameters and features

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

$b$  is a number

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

linear algebra: count from 1

**NumPy** 

```
w = np.array([1.0, 2.5, -3.3])
b = 4
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

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

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

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

```
f = 0
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
```



Gradient descent

 $\vec{w} = (w_1 \ w_2 \ \dots \ w_{16}) \quad \cancel{b} \text{ parameters}$ 

derivatives

 $\vec{d} = (d_1 \ d_2 \ \dots \ d_{16})$ 

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

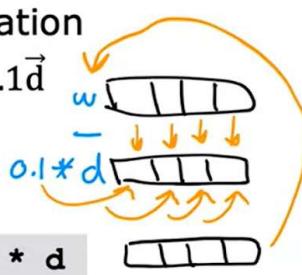
compute  $w_j = w_j - 0.1d_j$  for  $j = 1 \dots 16$

Without vectorization

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

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

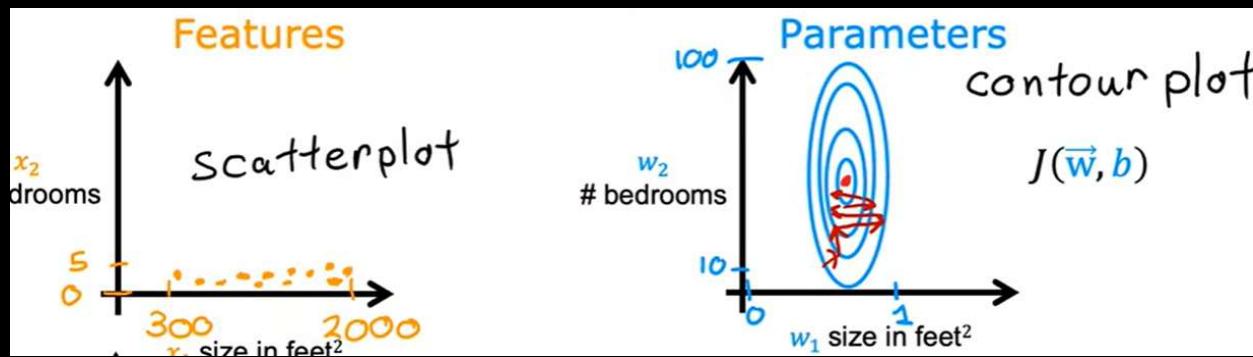
	Previous notation	Vector notation
Parameters	$w_1, \dots, w_n$ $b$	$\vec{w}$ vector of length n $b$ still a number
Model	$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$	$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
Cost function	$J(\underbrace{w_1, \dots, w_n}_J, b)$	$J(\underbrace{\vec{w}}_{\text{dot product}}, b)$
Gradient descent	<pre>repeat {     <math>w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\underbrace{w_1, \dots, w_n}_J, b)</math>     <math>b = b - \alpha \frac{\partial}{\partial b} J(\underbrace{w_1, \dots, w_n}_J, b)</math> }</pre>	
	<pre>repeat {     <math>w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)</math>     <math>b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)</math> }</pre>	

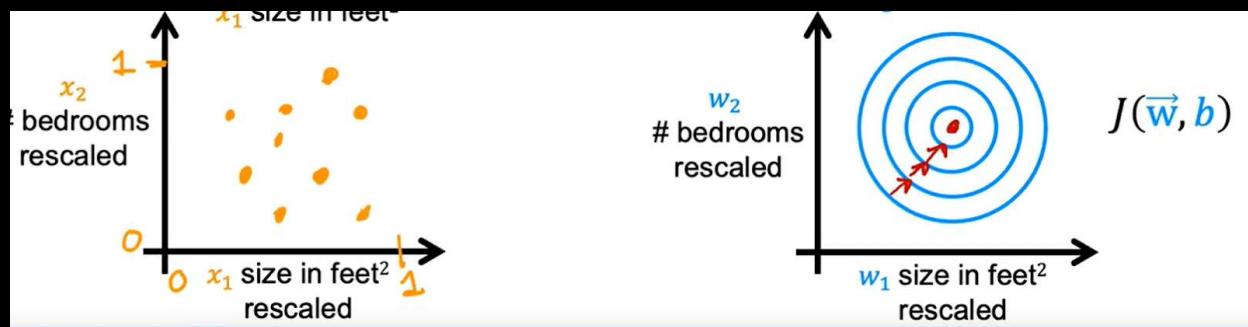
## Gradient descent

<p>repeat {</p> <p style="margin-left: 20px;">One feature</p> $\underline{w} = \underline{w} - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\underline{w}, b}(\underline{x}^{(i)}) - \underline{y}^{(i)}) \underline{x}^{(i)}$ <p style="margin-left: 40px;"><math>\hookrightarrow \frac{\partial}{\partial \underline{w}} J(\underline{w}, b)</math></p> $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\underline{w}, b}(\underline{x}^{(i)}) - \underline{y}^{(i)})$ <p style="margin-left: 20px;">simultaneously update <math>\underline{w}, b</math></p> <p>}</p>	<p><math>n</math> features (<math>n \geq 2</math>)</p> <p>repeat {</p> <p style="margin-left: 20px;"><math>j = 1</math></p> $\underline{w}_1 = \underline{w}_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - \underline{y}^{(i)}) \underline{x}_1^{(i)}$ <p style="margin-left: 40px;"><math>\hookrightarrow \frac{\partial}{\partial \underline{w}_1} J(\vec{w}, b)</math></p> <p style="margin-left: 20px;"><math>j = n</math></p> $\underline{w}_n = \underline{w}_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - \underline{y}^{(i)}) \underline{x}_n^{(i)}$ $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - \underline{y}^{(i)})$ <p style="margin-left: 20px;">simultaneously update <math>\underline{w}_j</math> (for <math>j = 1, \dots, n</math>) and <math>b</math></p> <p>}</p>
---	---

## Features and parameter values:

	size of feature $x_j$	size of parameter $w_j$
size in feet <sup>2</sup>	↔	↔
#bedrooms	↔	↔





## Feature scaling methods

$$300 \leq x_1 \leq 2000$$

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

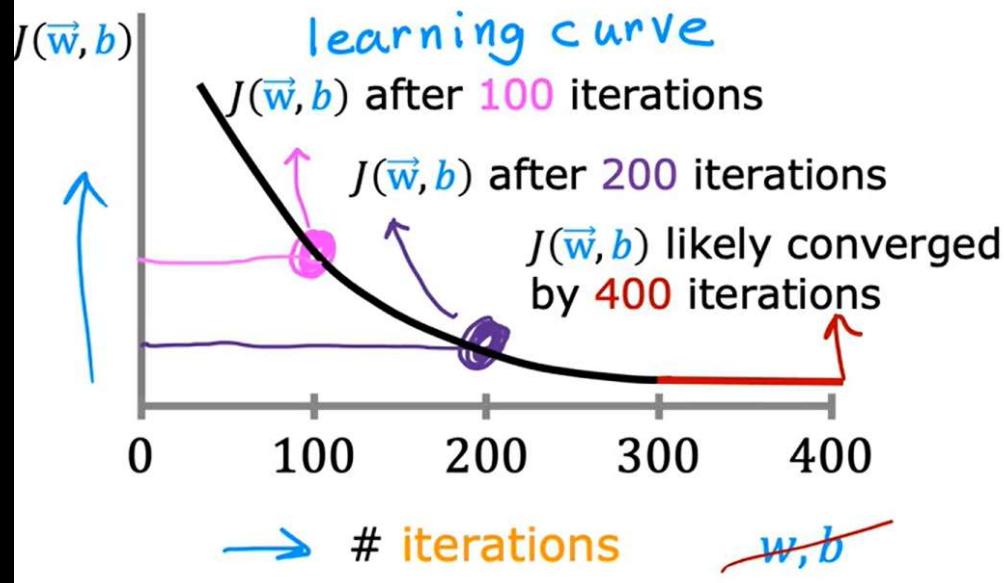
$$x_1 = \frac{x_1 - \mu_1}{\text{max-min}}$$

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

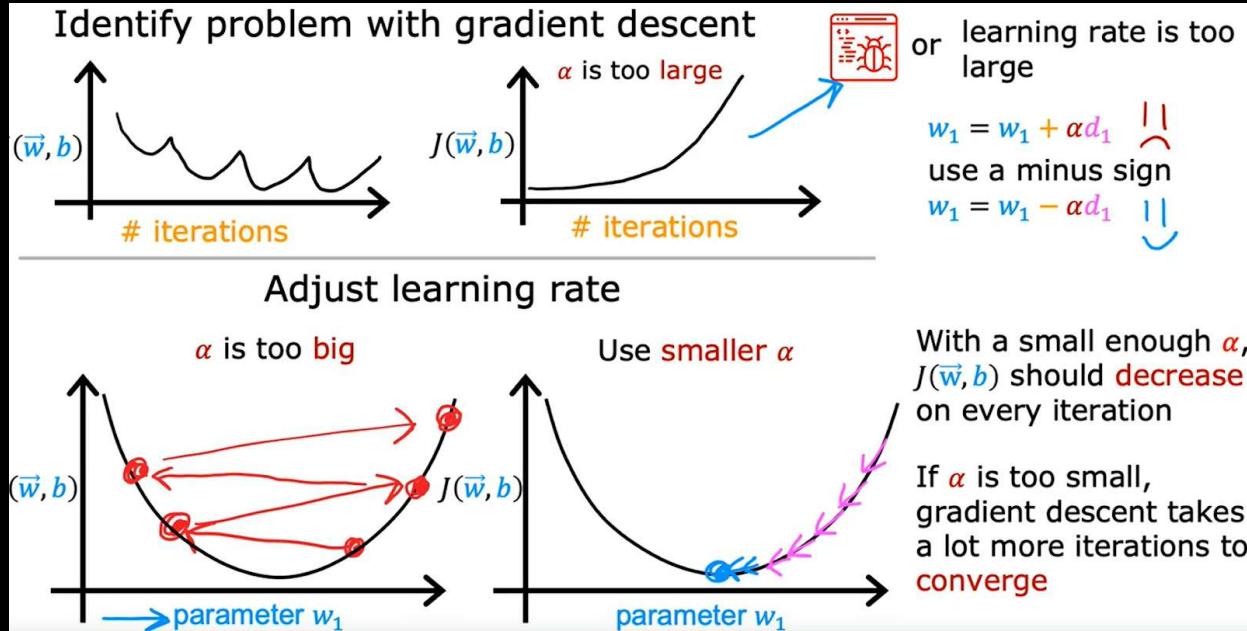
aim for about  $-1 \leq x_j \leq 1$  for each feature  $x_j$   
 $\quad -3 \leq x_j \leq 3$       } acceptable ranges  
 $\quad -0.3 \leq x_j \leq 0.3$

$0 \leq x_1 \leq 3$	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

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



# Learning rate(alpha)



... 0.001 0.003 0.01 0.03 0.1 0.3 1 ...  
   $\xrightarrow{3X}$     $\xrightarrow{\approx 3X}$     $\xrightarrow{3X}$     $\xrightarrow{\approx 3X}$     $\xrightarrow{3X}$     $\xrightarrow{\approx 3X}$

## Feature engineering

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

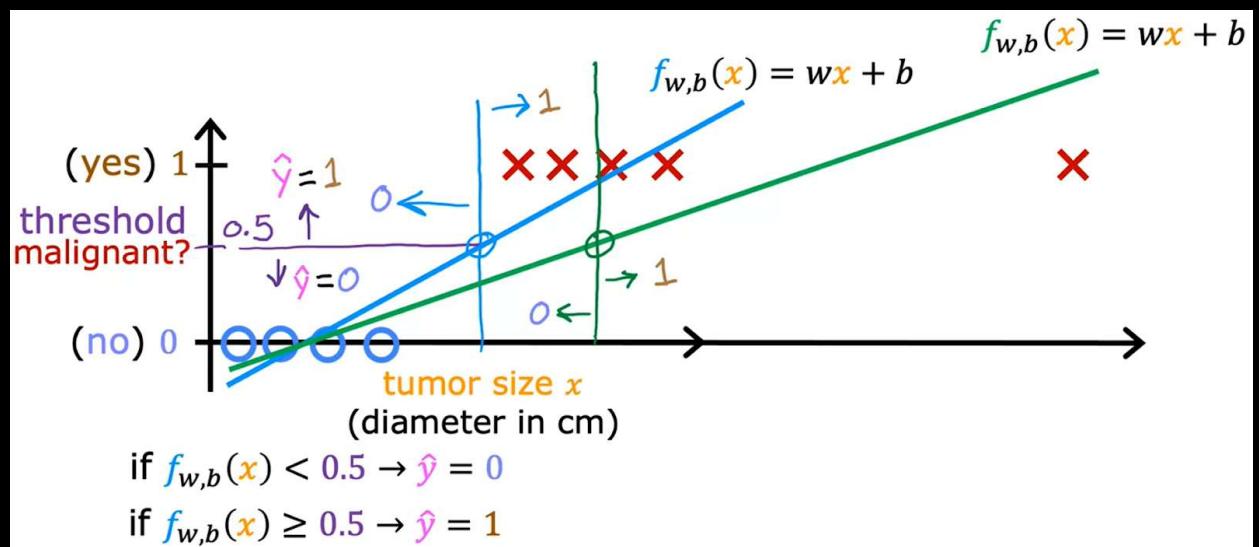
frontage      depth

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

$$x_3 = x_1 x_2$$

new feature

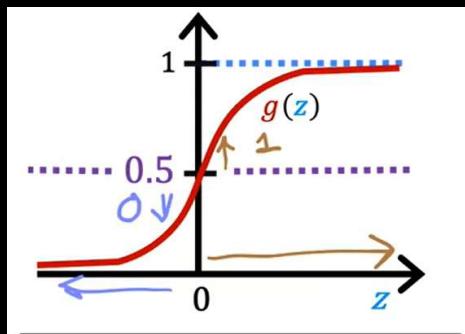
$$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$



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

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

"logistic regression"



$$f_{\vec{w}, b}(\vec{x}) = g\left(\frac{\vec{w} \cdot \vec{x} + b}{2}\right) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

$= P(Y=1|\vec{x}; \vec{w}, b)$  to decide 0 or 1

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

Yes:  $\hat{y} = 1$ . No:  $\hat{y} = 0$ . |  $\hat{y}$  = prediction

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

$$g(z) \geq 0.5$$

$$z \geq 0$$

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

$$\hat{y} = 1$$

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

$$\hat{y} = 0$$

### decision boundary

for ex. taking & plotting a graphs, that includes 2 features  $x_1$  &  $x_2$ :

$$\text{i.e., } f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1x_1 + w_2x_2 + b)$$

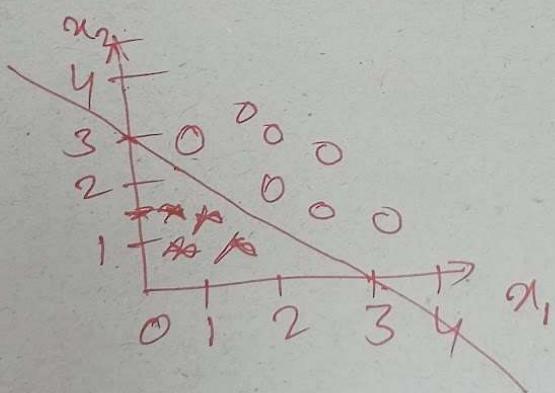
for ex. take  $w_1 = w_2 = 1$  &  $b = -3$ .

$$\text{decision boundary } z = \vec{w} \cdot \vec{x} + b = 0 \quad | \text{ equal to zero}$$

$$z = w_1 + w_2 + b = 0$$

$$w_1 + w_2 - 3 = 0$$

$$w_1 + w_2 = 3$$



all the data points that come left are into one class & other into another class

there are non linear decision boundaries are also possible, where  $z$  is complex. for ex

$$z = w_1x_1^2 + w_2x_2^2 + b. \text{ this is a circle.}$$

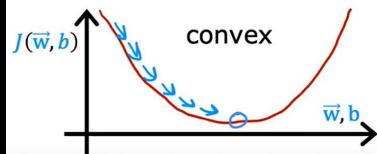
## Squared error cost

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

loss  $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$

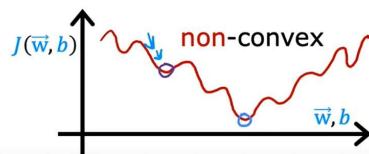
linear regression

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



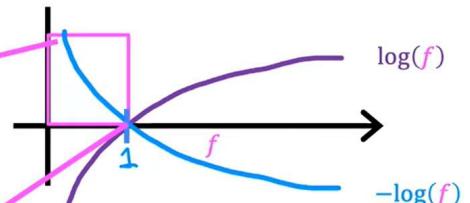
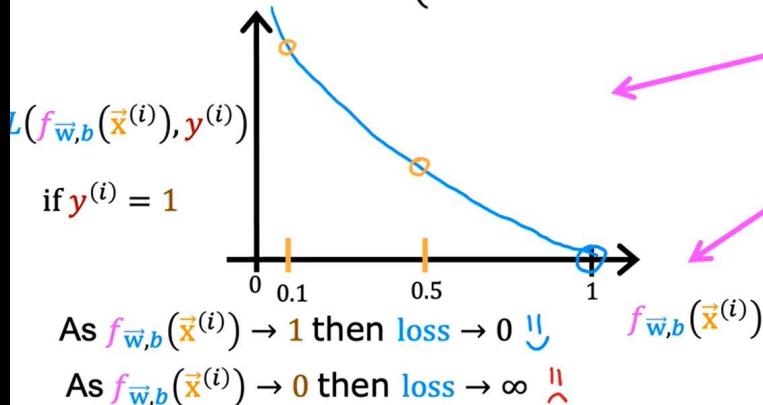
logistic regression

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

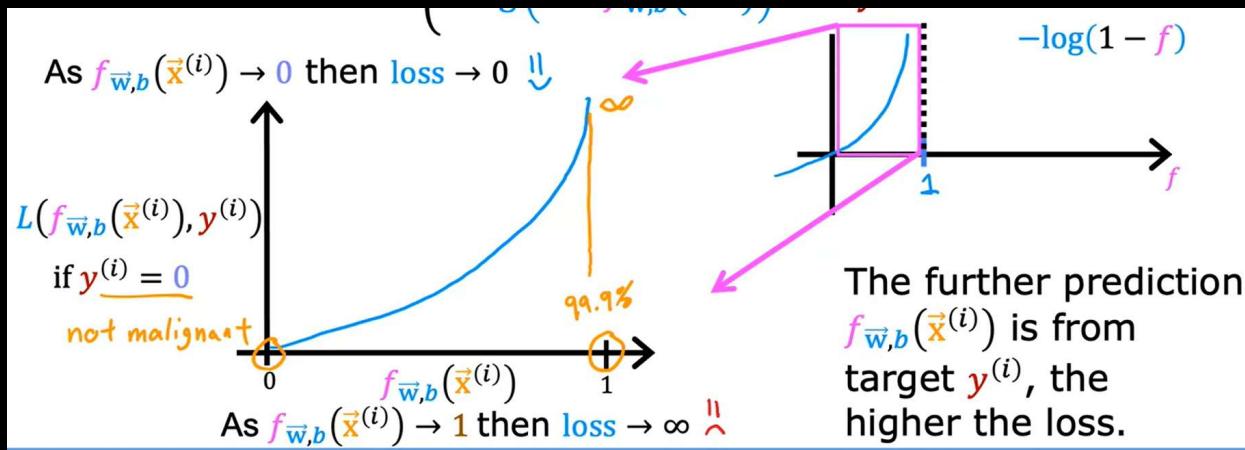


## Logistic loss function

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



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



$$J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \underbrace{L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})}_{\text{loss}}$$

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

convex  
can reach a global minimum

find  $w, b$  that minimize cost  $J$

## Simplified loss function

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

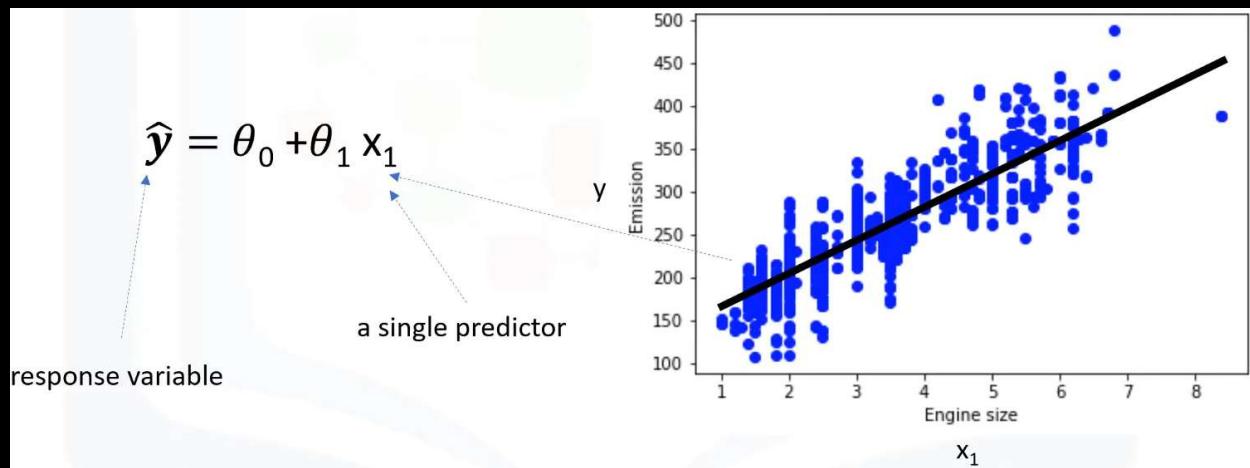
$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

if  $y^{(i)} = 1$ :  $\quad \quad \quad O \quad \quad \quad (1 - O)$

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

if  $y^{(i)} = 0$ :

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -\underbrace{(1 - O) \log(1 - f(\vec{x}))}_{}$$



$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

# Estimating the parameters

	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS
0	2.0	4	8.5	196
1	2.4	4	9.6	221
2	1.5	4	5.9	136
3	3.5	6	11.1	255
4	X <sub>1</sub>			y
5	3.5	6	10.6	244
6	3.5	6	10.0	230
7	3.5	6	10.1	232
8	3.7	6	11.1	255
	3.7	6	11.6	267

$$\hat{y} = \theta_0 + \theta_1 x_1$$

$$\theta_1 = \frac{\sum_{i=1}^s (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^s (x_i - \bar{x})^2}$$

$$\bar{x} = (2.0 + 2.4 + 1.5 + \dots) / 9 = 3.03$$

$$\bar{y} = (196 + 221 + 136 + \dots) / 9 = 226.22$$

$$\theta_1 = \frac{(2.0 - 3.03)(196 - 226.22) + (2.4 - 3.03)(221 - 226.22) + \dots}{(2.0 - 3.03)^2 + (2.4 - 3.03)^2 + \dots}$$

$$\theta_1 = 39$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$\theta_0 = 226.22 - 39 * 3.03$$

$$\theta_0 = 125.74$$

Model evaluation:

- The goal of regression is to build a model to accurately predict an unknown case.

train and test on same data set

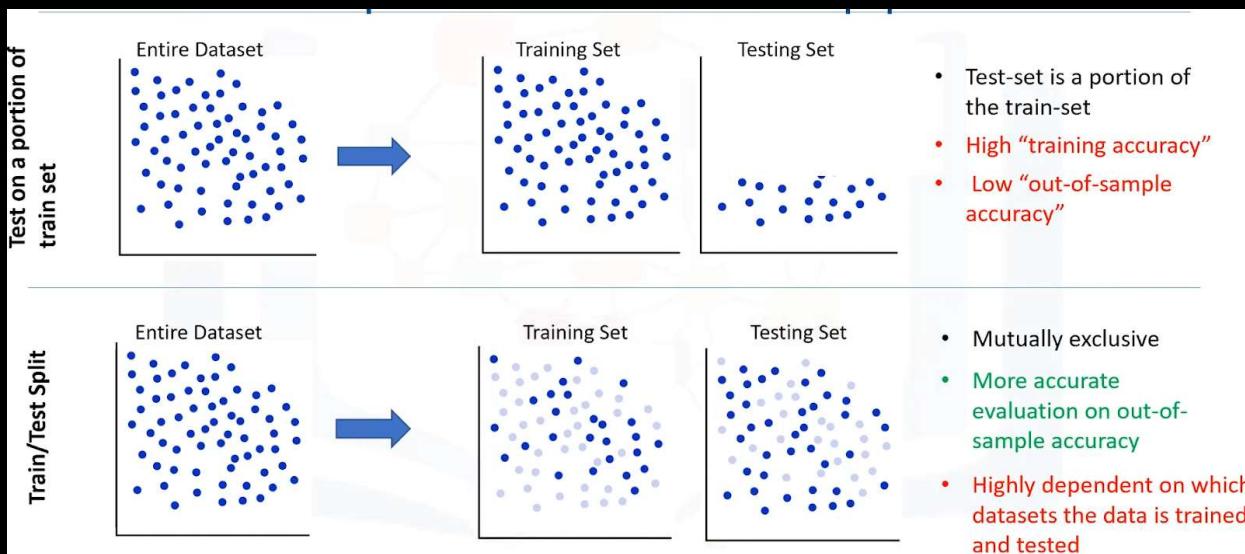
	ENGINESIZE	CYLINDERS	FUELCONSUMPTION_COMB	CO2EMISSIONS	Error = $\frac{(232 - 234) + (255 - 256) + \dots}{4}$
0	2.0	4	8.5	196	
1	2.4	4	9.6	221	
2	1.5	4	5.9	136	
3	3.5	6	11.1	255	
4	3.5	6	10.6	244	
5	3.5	6	10.0	230	
6	3.5	6	10.1	232	$Error = \frac{1}{n} \sum_{j=1}^n  y_j - \hat{y}_j $
7	3.7	6	11.1	255	
8	3.7	6	11.6	267	
9	2.4	4	9.2	212	

Test       $y$

Actual values

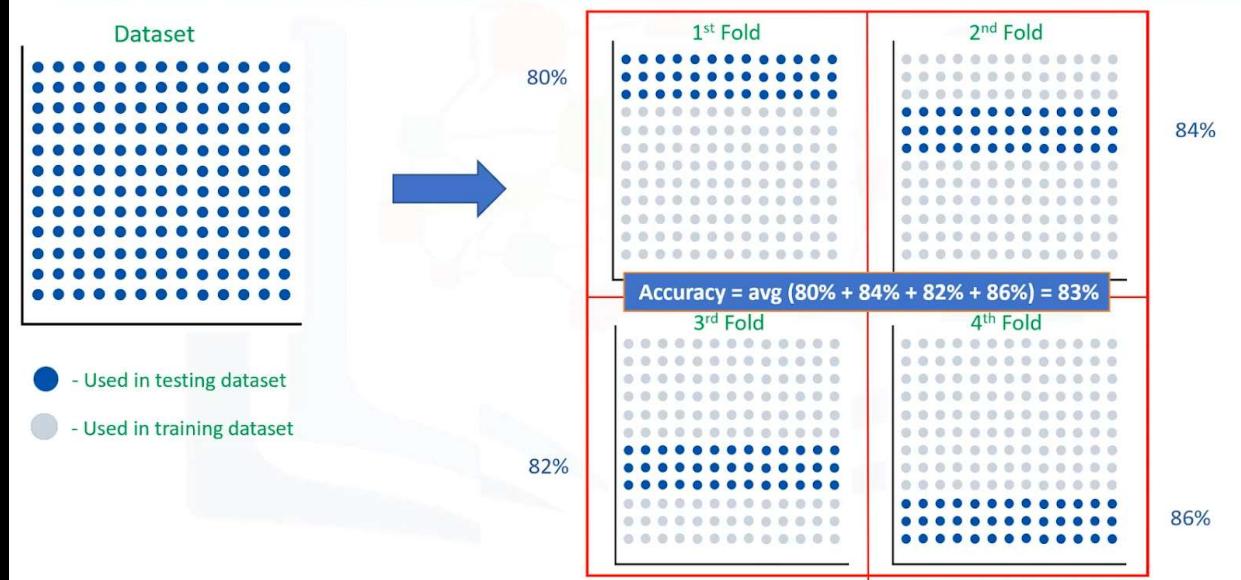
$\hat{y}$       Predicted values

**Train/test split** involves splitting the dataset into training and testing sets respectively, which are mutually exclusive.



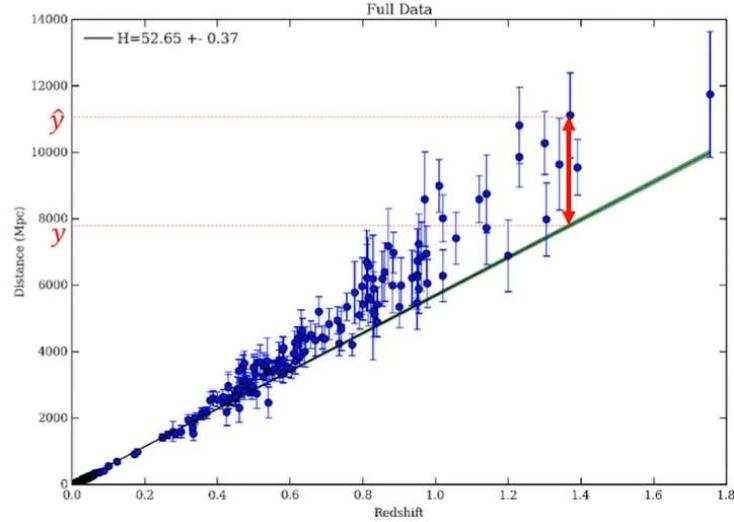
## K-fold cross-validation

# How to use K-fold cross-validation?



## Accuracy metrics for model evaluation

# What is an error of the model?



$$MAE = \frac{1}{n} \sum_{j=1}^n |y_j - \hat{y}_j|$$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

$$RAE = \frac{\sum_{j=1}^n |y_j - \hat{y}_j|}{\sum_{j=1}^n |y_j - \bar{y}|}$$

$$RSE = \frac{\sum_{j=1}^n (y_j - \hat{y}_j)^2}{\sum_{j=1}^n (y_j - \bar{y})^2}$$

$$R^2 = 1 - RSE$$

