

1. Classification motivation

## 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

false    true

0    1

useful for classification

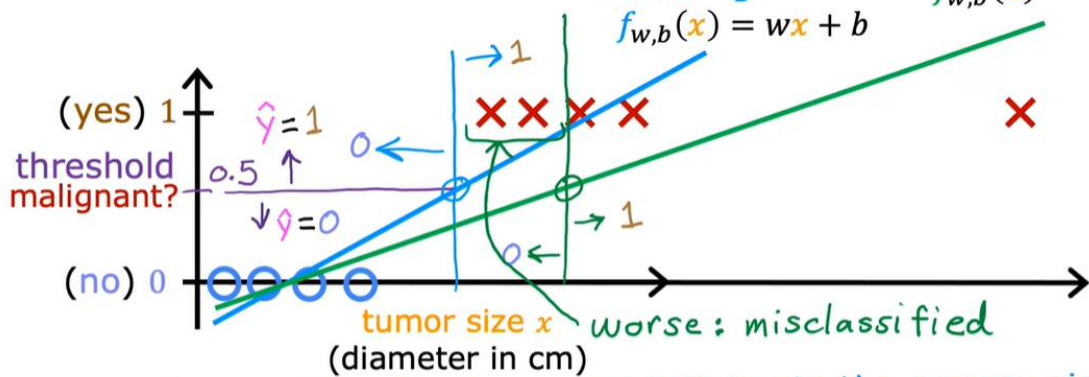
"negative class"  
≠ "bad"  
absence

"positive class"  
≠ "good"  
presence

decision boundary

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

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

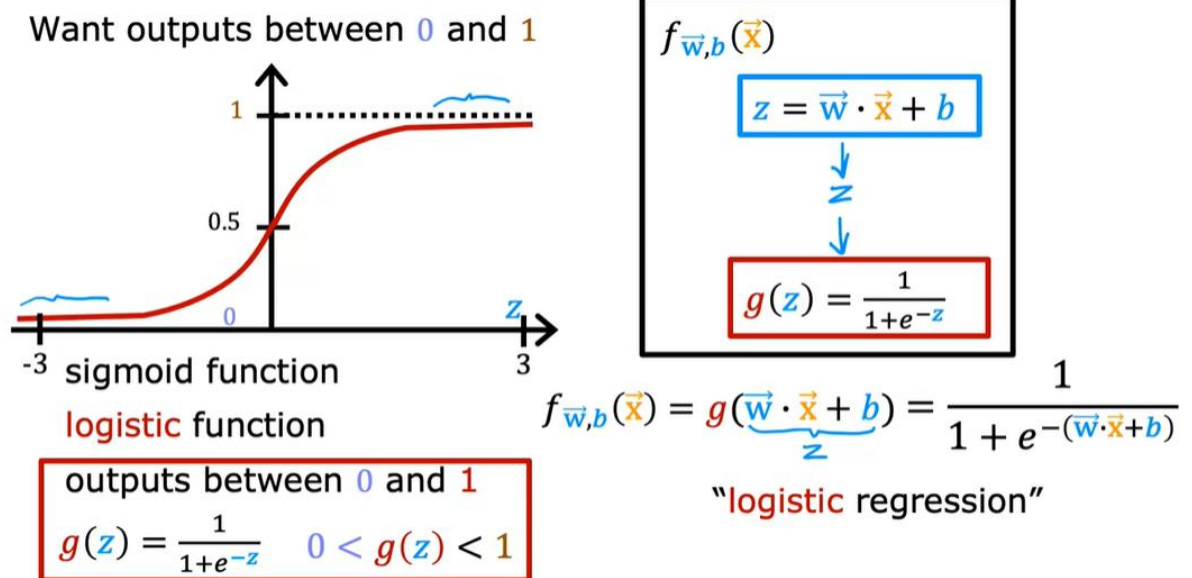
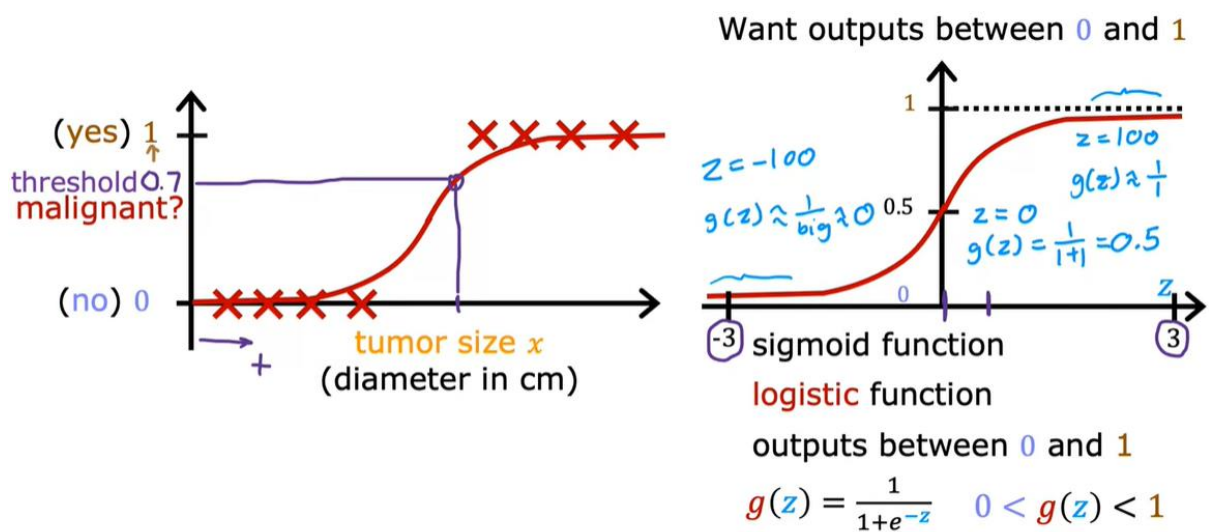


if  $f_{w,b}(x) < 0.5 \rightarrow \hat{y} = 0$

if  $f_{w,b}(x) \geq 0.5 \rightarrow \hat{y} = 1$

next: logistic regression

2. Logistic regression



## 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

Example:

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

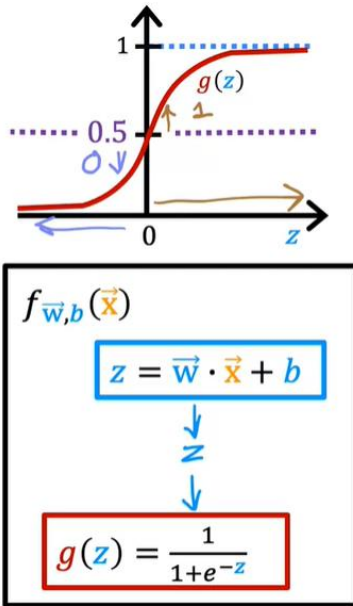
$f_{\vec{w},b}(\vec{x}) = 0.7$   
70% chance that  $y$  is 1

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

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

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

### 3. Decision boundary



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

$$= P(y = 1 | \vec{x}; \vec{w}, b) \quad 0.7 \quad 0.3$$

0 or 1? threshold

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

Yes:  $\hat{y} = 1$

No:  $\hat{y} = 0$

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

$$g(z) \geq 0.5$$

$$z \geq 0$$

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

$$\hat{y} = 1$$

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

$$\hat{y} = 0$$

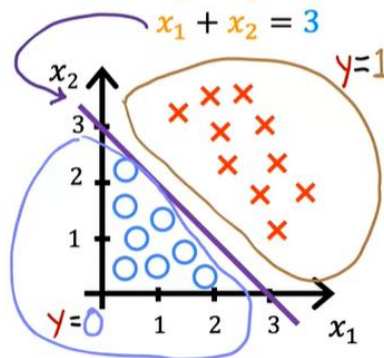
## Decision boundary

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

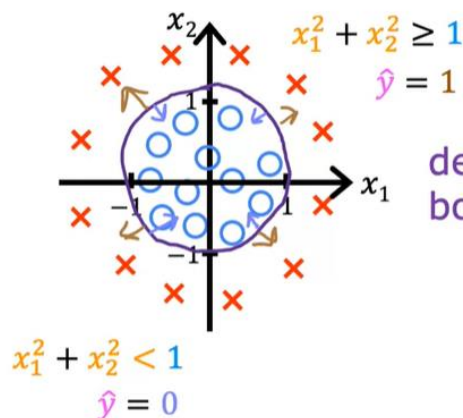
Decision boundary  $z = \vec{w} \cdot \vec{x} + b = 0$

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

$$x_1 + x_2 = 3$$



## Non-linear decision boundaries

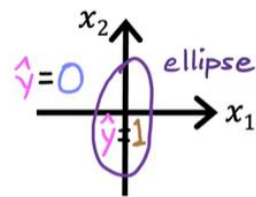


$$f_{\vec{w},b}(\vec{x}) = g(z) = g(\underbrace{w_1 x_1^2 + w_2 x_2^2 + b}_z)$$

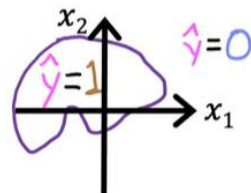
decision boundary  $z = x_1^2 + x_2^2 - 1 = 0$

boundary  $x_1^2 + x_2^2 = 1$

# Non-linear decision boundaries



$$f_{\vec{w},b}(\vec{x}) = g(z) = g(w_1x_1 + w_2x_2 + w_3x_1^2 + w_4x_1x_2 + w_5x_2^2 + w_6x_1^3 + \dots + b)$$



Let's say you are creating a tumor detection algorithm. Your algorithm will be used to flag potential tumors for future inspection by a specialist. What value should you use for a threshold?

- ☐ High, say a threshold of 0.9?
- ☒ Low, say a threshold of 0.2?

✓ **Correct**

**Correct:** You would not want to miss a potential tumor, so you will want a low threshold. A specialist will review the output of the algorithm which reduces the possibility of a 'false positive'. The key point of this question is to note that the threshold value does not need to be 0.5.

## 4. Practice quiz 1

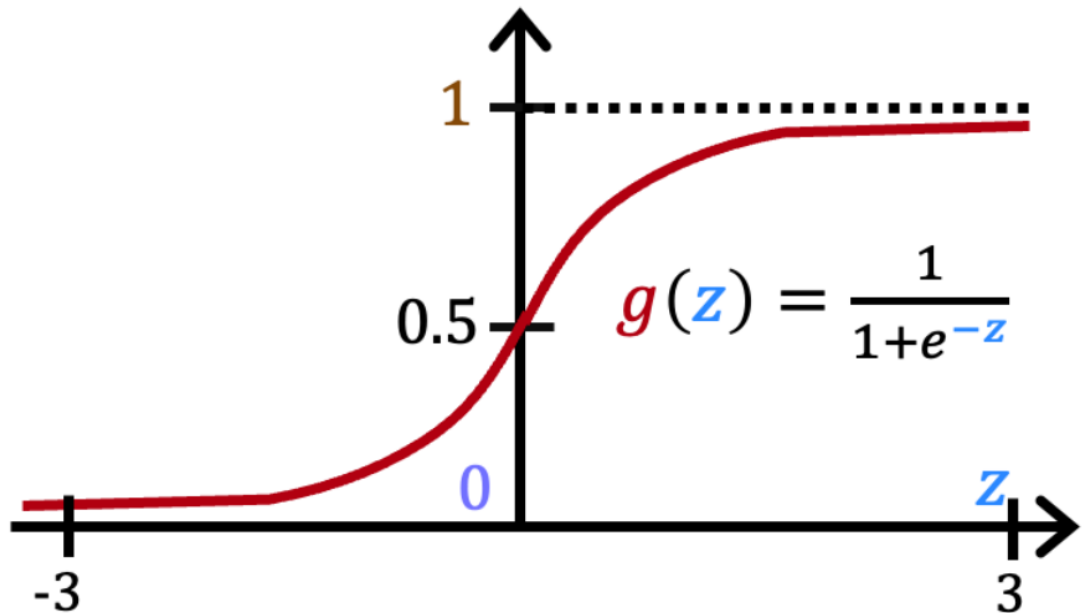
1. Which is an example of a classification task?
  - ☒ Based on the size of each tumor, determine if each tumor is malignant (cancerous) or not.
  - ☐ Based on a patient's blood pressure, determine how much blood pressure medication (a dosage measured in milligrams) the patient should be prescribed.
  - ☐ Based on a patient's age and blood pressure, determine how much blood pressure medication (measured in milligrams) the patient should be prescribed.

✓ **Correct**

This task predicts one of two classes, malignant or not malignant.

2. Recall the sigmoid function is  $g(z) = \frac{1}{1+e^{-z}}$

## sigmoid function



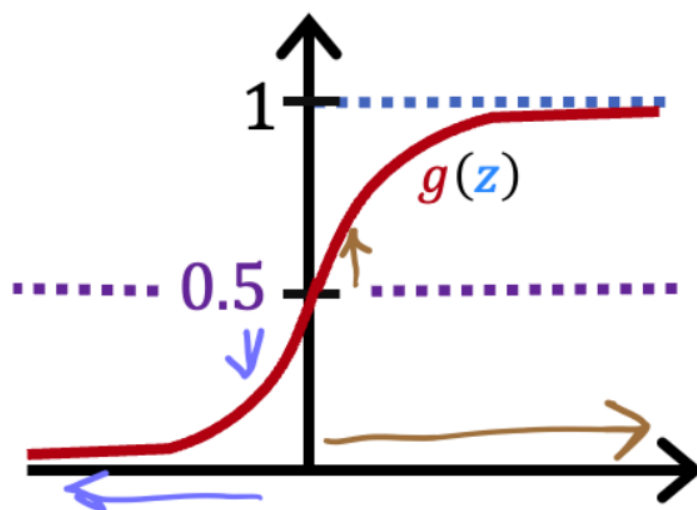
If  $z$  is a large positive number, then:

- ☐  $g(z)$  will be near 0.5
- ☐  $g(z)$  is near negative one (-1)
- ☐  $g(z)$  will be near zero (0)
- ☒  $g(z)$  is near one (1)

✓ Correct

Say  $z = +100$ . So  $e^{-z}$  is then  $e^{-100}$ , a really small positive number. So,  $g(z) = \frac{1}{1+\text{a small positive number}}$  which is close to 1

3.



A cat photo classification model predicts 1 if it's a cat, and 0 if it's not a cat. For a particular photograph, the logistic regression model outputs  $g(z)$  (a number between 0 and 1). Which of these would be a reasonable criteria to decide whether to predict if it's a cat?

- ☐ Predict it is a cat if  $g(z) < 0.5$
- ☐ Predict it is a cat if  $g(z) = 0.5$
- ☒ Predict it is a cat if  $g(z) \geq 0.5$
- ☐ Predict it is a cat if  $g(z) < 0.7$

✓ Correct

Think of  $g(z)$  as the probability that the photo is of a cat. When this number is at or above the threshold of 0.5, predict that it is a cat.

4.

True/False? No matter what features you use (including if you use polynomial features), the decision boundary learned by logistic regression will be a linear decision boundary.

- ☒ False
- ☐ True

✓ Correct

The decision boundary can also be non-linear, as described in the lectures.

5. Cost function for logistic regression

## Training set

	tumor size (cm)	...	patient's age	malignant?	$i = 1, \dots, m \leftarrow \text{training examples}$
	$x_1$		$x_n$	$y$	$j = 1, \dots, n \leftarrow \text{features}$
$i=1$	10		52	1	<div style="border: 1px solid red; padding: 2px; display: inline-block;">target <math>y</math> is 0 or 1</div> $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$
$\vdots$	2		73	0	
$\vdots$	5		55	0	
$i=m$	12		49	1	
	...		...	...	

How to choose  $\vec{w} = [w_1 \ w_2 \ \dots \ w_n]$  and  $b$ ?



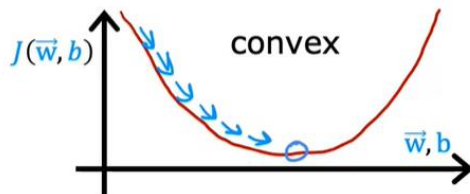
## Squared error cost

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

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

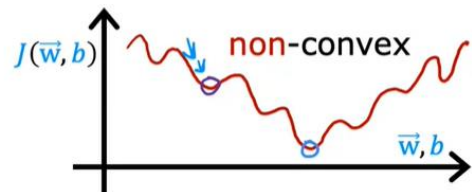
linear regression

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



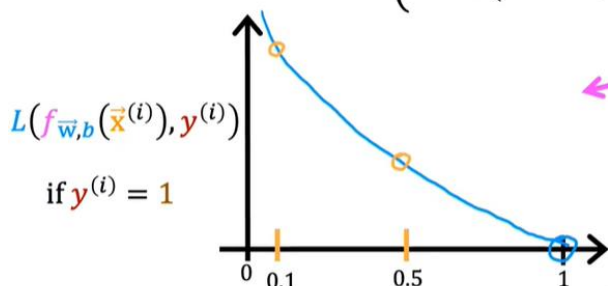
logistic regression

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



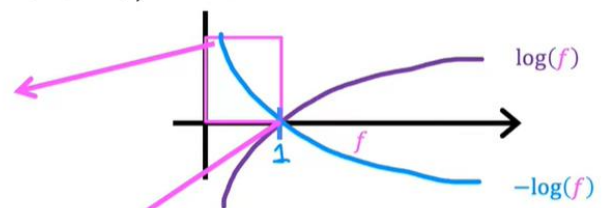
## Logistic loss function

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



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

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

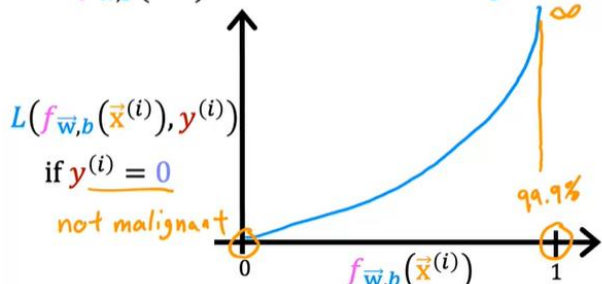


Loss is lowest when  $f_{\vec{w}, b}(\vec{x}^{(i)})$  predicts close to true label  $y^{(i)}$ .

## Logistic loss function

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

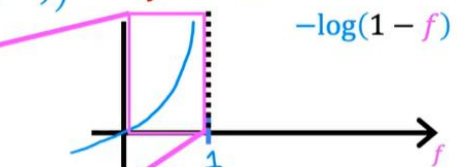
As  $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 0$  then loss  $\rightarrow 0$   $\Downarrow$



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

not malignant

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



The further prediction  $f_{\vec{w}, b}(\vec{x}^{(i)})$  is from target  $y^{(i)}$ , the higher the loss.

# Cost

cost

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

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

if  $y^{(i)} = 1$  convex  $\rightarrow$  can reach a global minimum  
if  $y^{(i)} = 0$  global minimum

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

Why is the squared error cost not used in logistic regression?

- ☐ The non-linear nature of the model results in a "wiggly", non-convex cost function with many potential local minima.
- ☐ The mean squared error is used for logistic regression.

☒ Correct

If using the mean squared error for logistic regression, the cost function is "non-convex", so it's more difficult for gradient descent to find an optimal value for the parameters  $w$  and  $b$ .

## 6. Simplified cost function for logistic regression

### Simplified loss function

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

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = - \underbrace{y^{(i)}}_0 \log(\underbrace{f_{\vec{w}, b}(\vec{x}^{(i)})}_{(1-0)}) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

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

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

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

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = - (1-0) \log(1 - f(\hat{x}))$$



## Simplified cost function

$$\begin{aligned}
 \text{loss} \\
 L(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)}) &= -y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) \\
 \text{cost} \\
 J(\vec{w}, b) &= \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w},b}(\vec{x}^{(i)}), y^{(i)})] \quad \text{convex (single global minimum)} \\
 &= \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w},b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w},b}(\vec{x}^{(i)}))] \\
 &\quad \text{maximum likelihood (don't worry about it!)}
 \end{aligned}$$

For the simplified loss function:

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

if the target  $y^{(i)} = 1$ , then what does this expression simplify to?

- ☐  $-\log(1 - f_{\vec{w},b}(\mathbf{x}^{(i)}))$   
☒  $-\log(f_{\vec{w},b}(\mathbf{x}^{(i)}))$

✓ Correct

The second term of the expression is reduced to zero when the target equals 1.

### 7. Practice quiz 2

1.

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

In this lecture series, "cost" and "loss" have distinct meanings. Which one applies to a single training example?

☒ Loss

✓ Correct

In these lectures, loss is calculated on a single training example. It is worth noting that this definition is not universal. Other lecture series may have a different definition.

☐ Cost

☐ Both Loss and Cost

☐ Neither Loss nor Cost

2.

Simplified loss function

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

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

For the simplified loss function, if the label  $y^{(i)} = 0$ , then what does this expression simplify to?

- ☒  $-\log(1 - f_{\vec{w},b}(\vec{x}^{(i)}))$
- ☐  $\log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) + \log(1 - f_{\vec{w},b}(\vec{x}^{(i)}))$
- ☐  $-\log(1 - f_{\vec{w},b}(\vec{x}^{(i)})) - \log(1 - f_{\vec{w},b}(\vec{x}^{(i)}))$
- ☐  $\log(f_{\vec{w},b}(\vec{x}^{(i)}))$

✓ **Correct**  
When  $y^{(i)} = 0$ , the first term reduces to zero.

8. Gradient descent implementation

## Training logistic regression

Find  $\vec{w}, b$

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

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

## 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 ... n*

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

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

} simultaneous updates

## 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:

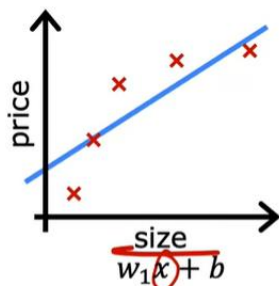
- Monitor gradient descent (learning curve)
- Vectorized implementation
- Feature scaling

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

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

### 9. The problem of overfitting

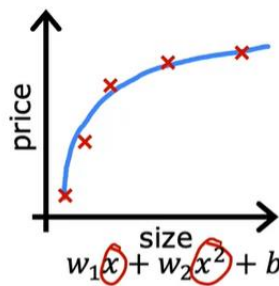
#### Regression example



*underfit*

- Does not fit the training set well

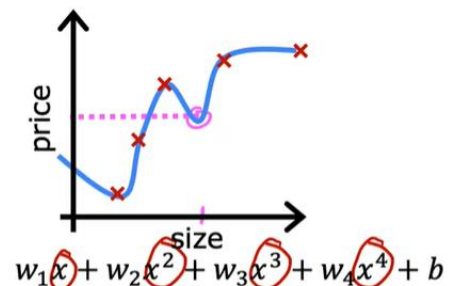
*high bias*



*just right*

- Fits training set pretty well

*generalization*

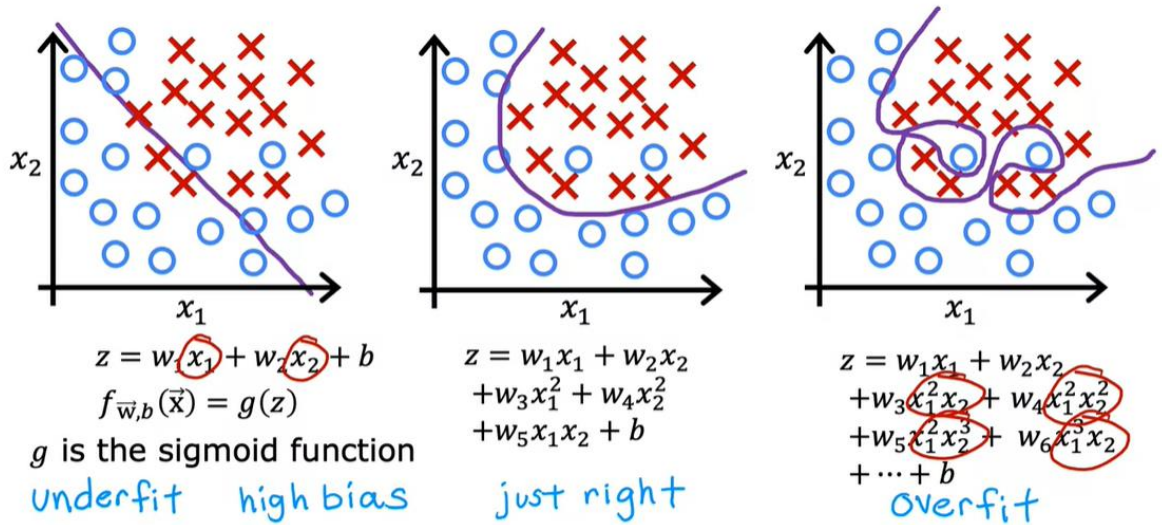


*overfit*

- Fits the training set extremely well

*high variance*

# Classification



Our goal when creating a model is to be able to use the model to predict outcomes correctly for **new examples**. A model which does this is said to **generalize** well.

When a model fits the training data well but does not work well with new examples that are not in the training set, this is an example of:

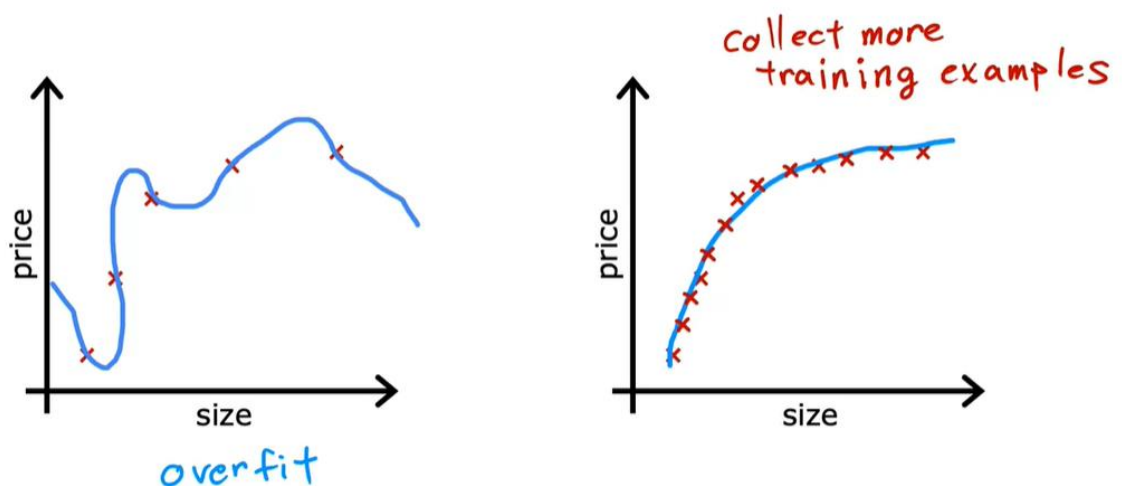
- ☐ A model that generalizes well (neither high variance nor high bias)
- ☐ Underfitting (high bias)
- ☒ Overfitting (high variance)
- ☐ None of the above

✓ Correct

This is when the model does not generalize well to new examples.

## 10. Addressing overfitting

### Collect more training examples

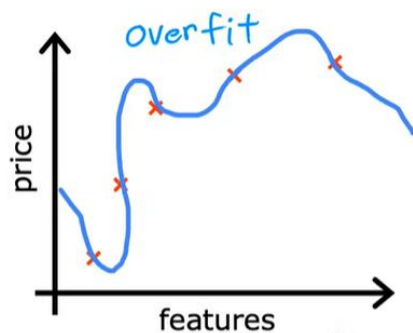


## Select features to include/exclude



## Regularization

Reduce the size of parameters  $w_j$



$$f(x) = 28x - 385x^2 + 39x^3 - 174x^4 + 100$$

$\underbrace{0}$  ← eliminate feature  
 large values for  $w_j$



$$f(x) = 13x - 0.23x^2 + 0.000014x^3 - 0.0001x^4 + 10$$

small values for  $w_j$

## Addressing overfitting

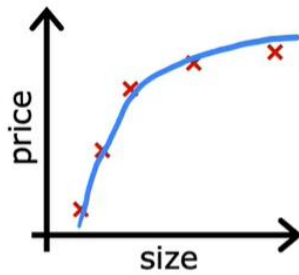
### Options

1. Collect more data
2. Select features
  - Feature selection in course 2
3. Reduce size of parameters
  - "Regularization" next videos!

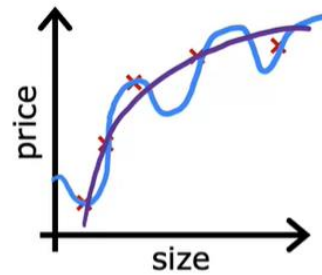


## 11. Cost function with regularization

### Intuition



$$w_1x + w_2x^2 + b$$



$$w_1x + w_2x^2 + \underbrace{w_3x^3}_{\approx 0} + \underbrace{w_4x^4}_{\approx 0} + b$$

make  $w_3, w_4$  really small ( $\approx 0$ )

$$\min_{\vec{w}, b} \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \cancel{1000 \underbrace{w_3^2}_{0.001}} + \cancel{1000 \underbrace{w_4^2}_{0.002}}$$

### Regularization

small values  $w_1, w_2, \dots, w_n, b$

simpler model

less likely to overfit

$$w_3 \approx 0$$

$$w_4 \approx 0$$

size	bedrooms	floors	age	avg income	...	distance to coffee shop	price
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$		$x_{100}$	$y$
$w_1, w_1, w_2, \dots, w_{100}, b$							
$n$ features							
$n = 100$							

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

"lambda" regularization parameter  $\lambda > 0$

can include or exclude  $b$



# Regularization

mean squared error      regularization term

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

$\lambda$  balances both goals

choose  $\lambda = 10^{10}$

$$f_{\vec{w}, b}(\vec{x}) = \cancel{w_1 x} + \cancel{w_2 x^2} + \cancel{w_3 x^3} + \cancel{w_4 x^4} + b$$

$\approx 0$        $\approx 0$        $\approx 0$        $\approx 0$

$f(x) = b$       choose  $\lambda$

For a model that includes the regularization parameter  $\lambda$  (lambda), increasing  $\lambda$  will tend to...

- ☐ Increase the size of parameter  $b$ .
- ☐ Decrease the size of the parameter  $b$ .
- ☒ Decrease the size of parameters  $w_1, w_2, \dots, w_n$ .
- ☐ Increases the size of the parameters  $w_1, w_2, \dots, w_n$ .

✓ Correct

Increasing the regularization parameter *lambda* reduces overfitting by reducing the size of the parameters. For some parameters that are near zero, this reduces the effect of the associated features.

## 12. Regularized linear regression

### 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  $b$

How we get the derivative term (optional)

$$\begin{aligned}
 \frac{\partial}{\partial w_j} J(\vec{w}, b) &= \frac{\partial}{\partial w_j} \left[ \frac{1}{2m} \sum_{i=1}^m \underbrace{(f(\vec{x}^{(i)}) - y^{(i)})^2}_{\vec{w} \cdot \vec{x}^{(i)} + b} + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right] \\
 &= \frac{1}{2m} \sum_{i=1}^m \left[ (\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}) \cancel{2x_j^{(i)}} \right] + \frac{\lambda}{2m} \cancel{2} w_j \quad \text{No } \sum_{j=1}^n \\
 &= \frac{1}{m} \sum_{i=1}^m \left[ \underbrace{(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)})}_{f(\vec{x})} x_j^{(i)} \right] + \frac{\lambda}{m} w_j \\
 &= \frac{1}{m} \sum_{i=1}^m \left[ (f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j
 \end{aligned}$$

Implementing gradient descent

repeat {

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

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

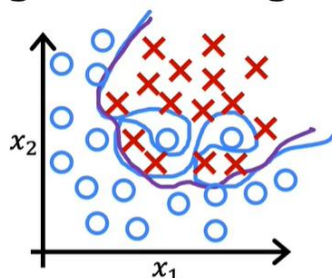
} simultaneous update  $j = 1 \dots n$

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

$$\begin{aligned}
 \alpha \frac{\lambda}{m} &= 0.01 \frac{1}{50} = 0.0002 \\
 w_j (1 - 0.0002) &= 0.9998 w_j
 \end{aligned}$$

### 13. Regularized logistic regression

Regularized logistic regression



$$\begin{aligned}
 z &= w_1 x_1 + w_2 x_2 \\
 &\quad + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 \\
 &\quad + w_5 x_1^2 x_2^3 + \dots + b
 \end{aligned}$$

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

Cost function

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

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

# Regularized logistic regression

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

*min  $\vec{w}, b$*

## Gradient descent

repeat {

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

*j = 1...n*

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

}

*Looks same as for linear regression!*

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

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

*logistic regression*

*don't have to regularize*

## 14. Practice quiz

1. Which of the following can address overfitting?

☒ Select a subset of the more relevant features.

☒ **Correct**

If the model trains on the more relevant features, and not on the less useful features, it may generalize better to new examples.

☒ Collect more training data

☒ **Correct**

If the model trains on more data, it may generalize better to new examples.

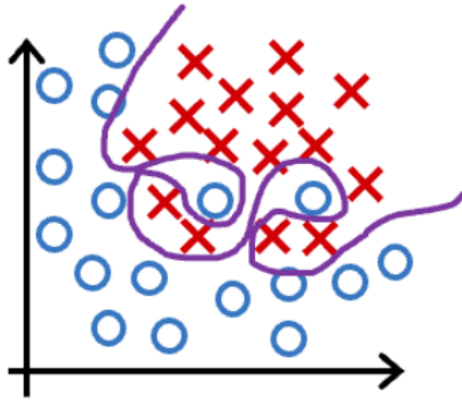
☒ Apply regularization

☒ **Correct**

Regularization is used to reduce overfitting.

☐ Remove a random set of training examples

2. You fit logistic regression with polynomial features to a dataset, and your model looks like this.



What would you conclude? (Pick one)

- ☐ The model has high bias (underfit). Thus, adding data is likely to help
- ☐ The model has high variance (overfit). Thus, adding data is, by itself, unlikely to help much.
- ☒ The model has high variance (overfit). Thus, adding data is likely to help
- ☐ The model has high bias (underfit). Thus, adding data is, by itself, unlikely to help much.

✓ Correct

The model has high variance (it overfits the training data). Adding data (more training examples) can help.

### 3. Regularization

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

Suppose you have a regularized linear regression model. If you increase the regularization parameter  $\lambda$ , what do you expect to happen to the parameters  $w_1, w_2, \dots, w_n$ ?

- ☐ This will increase the size of the parameters  $w_1, w_2, \dots, w_n$
- ☒ This will reduce the size of the parameters  $w_1, w_2, \dots, w_n$

✓ Correct

Regularization reduces overfitting by reducing the size of the parameters  $w_1, w_2, \dots, w_n$ .