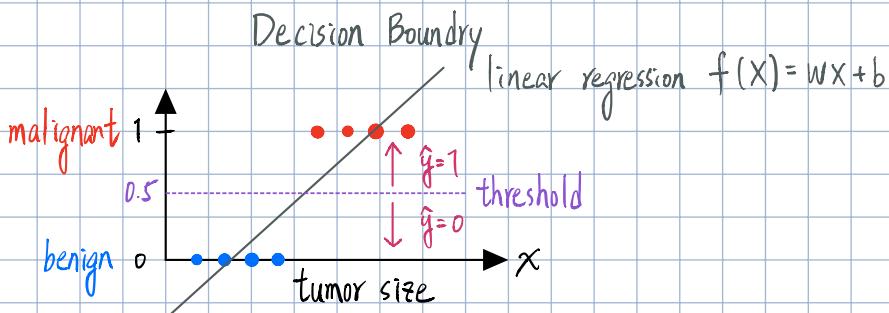
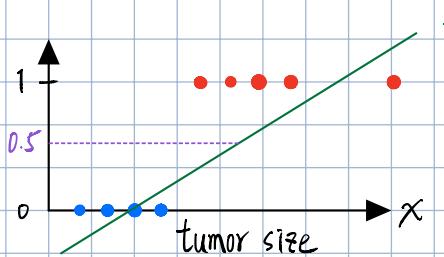


# Classification with Logistic Regression

Motivation: binary classification



Using threshold here makes sense. But, there's a new malignant sample as shown below.

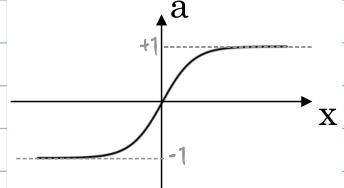


The new line destroys the threshold approach, making it hard to determine the threshold value.

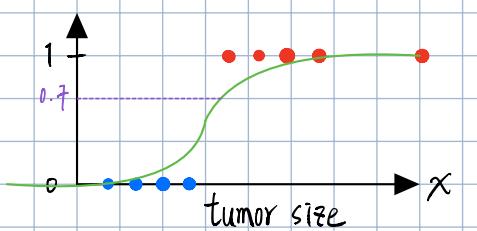
This is why linear regression doesn't work on classification problems. We need **Logistic Regression**.

Hyperbolic tangent:

$$a = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

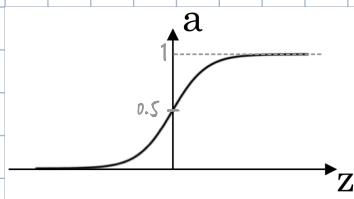


Logistic Regression



$$f(x) = g(w \cdot x + b) = g(z) = \frac{1}{1 + e^{-z}}$$

Sigmoid/Logistic function:



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

Interpretation of logistic regression output:

$$f_{w,b}(x) = g(w \cdot x + b) \rightarrow \text{Probability that class is 1. } P(y=0) + P(y=1) = 1$$

E.g.:  $x$  is tumor size  
 $y$  is 0 benign  
is 1 malignant

Notation:  $f(x) = P(y=1 | x; w, b)$

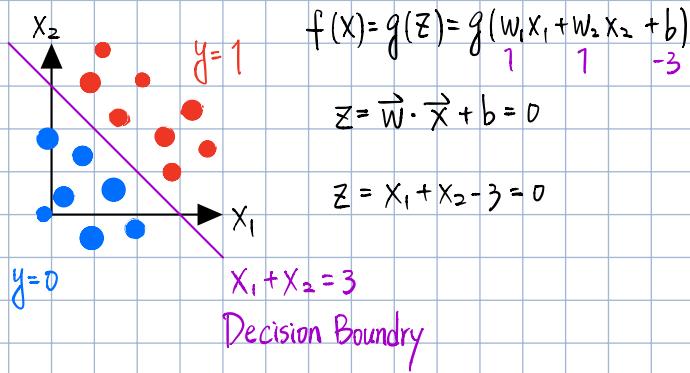
The probability that  $y$  is 1 given  $x$ , parameterized by  $w, b$ .

$f(x) = 0.7$ . 70% chance that  $y$  is 1.

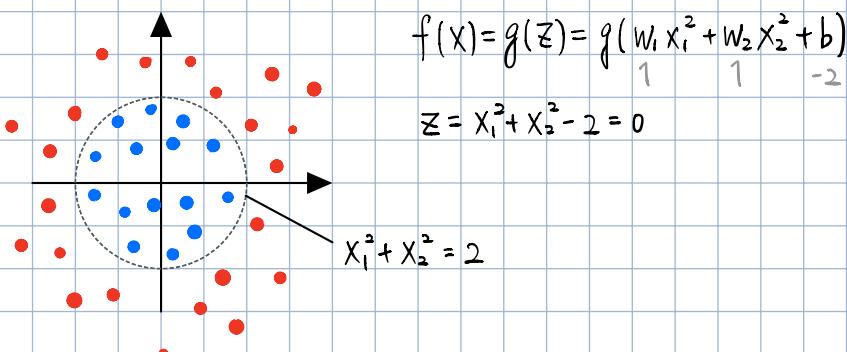
## Decision Boundary

$$\text{Given } f(X) = g(Z) = \frac{1}{1+e^{-Z}} = P(Y=1|X; w, b)$$

$\hat{y} = 1 \rightarrow f(X) \geq 0.5 \rightarrow g(Z) \geq 0.5 \rightarrow Z \geq 0 \rightarrow \boxed{wX + b \geq 0}$  Decision Boundary



nonlinear decision boundary



## Cost Function for Logistic Regression

Recall Squared Error Cost

Given  $m$  training samples ( $i=1, \dots, m$ ), and each sample has  $n$  features ( $j=1, \dots, n$ )

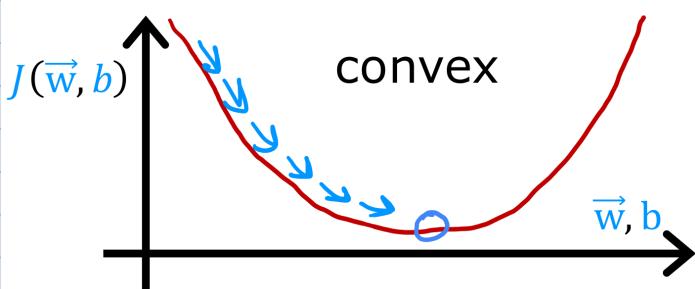
target  $y$  is 0 or 1.  $f(\vec{x}^{(i)}) = g(z) = \frac{1}{1+e^z}$ .  $0 < f(\vec{x}^{(i)}) < 1$

$$J(W, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f(\vec{x}^{(i)}) - y^{(i)})^2 \rightarrow \text{loss } L(f(\vec{x}^{(i)}), y^{(i)})$$

The two images below show that why squared error cost function is not a good choice for logistic regression.

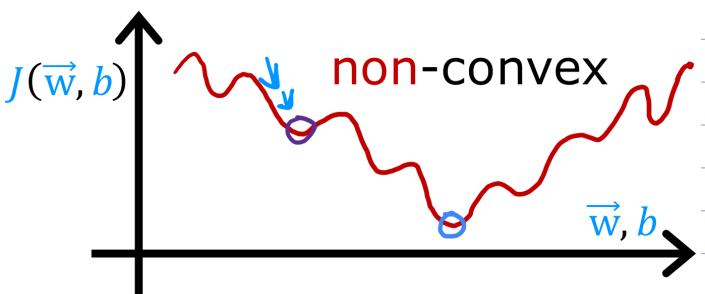
linear regression

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



logistic regression

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



Instead, a good choice of cost function for logistic regression is:

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

Try to explain why the loss makes sense here.

The loss function can be simplified as follow:

$$L(f_{w,b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f(\vec{x}^{(i)})) - (1-y^{(i)}) \log(1-f(\vec{x}^{(i)})) \quad \text{Convex, single global minimum}$$

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

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

# Gradient Descent for Logistic Regression

Find  $W, b$ .

Given new  $\vec{x}$ , output  $f(\vec{x}) = \frac{1}{1 + e^{-(wx+b)}}$ .  $P(y=1|X; W, b)$

Gradient Descent

$$\text{cost } J(W, b) = \frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f(x^{(i)})) + (1 - y^{(i)}) \log(1 - f(x^{(i)}))]$$

repeat until convergence {

$j = 1, \dots, n$

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

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

}

repeat until convergence {

$j = 1, \dots, n$

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

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

}

As a sidenote, in programming,

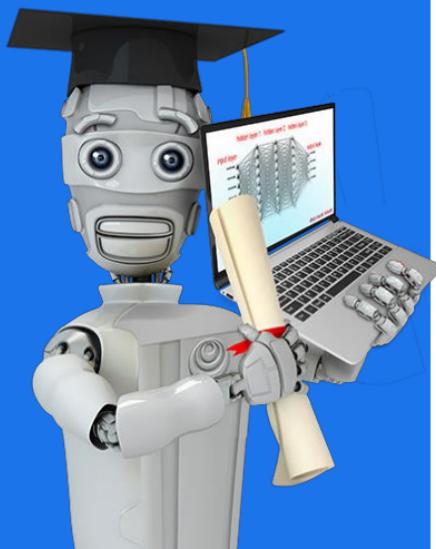
The name of the cost function for classification with exactly two classes is usually refer as

Binary Cross Entropy.

The name of the cost function for regression is usually refer as Mean Squared Error.

Stanford  
ONLINE

DeepLearning.AI

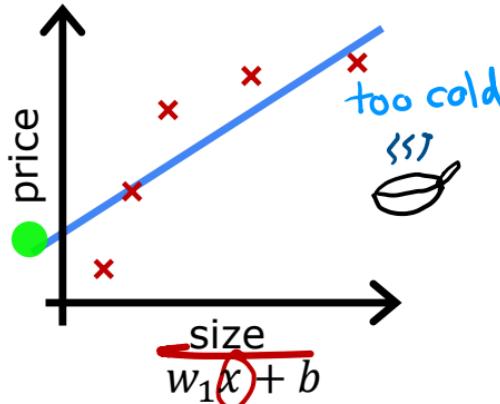


# Regularization to Reduce Overfitting

---

## The Problem of Overfitting

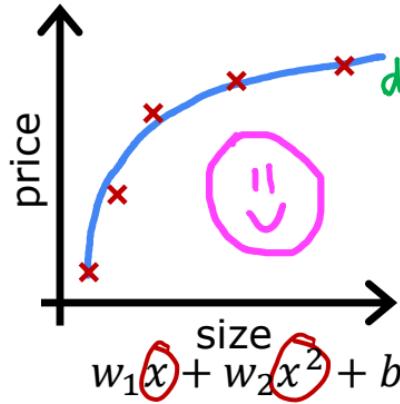
# Regression example



underfit

- Does not fit the training set well

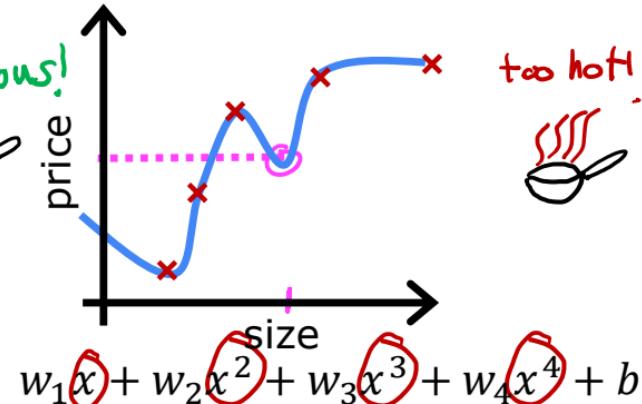
high bias



just right

- Fits training set pretty well

generalization

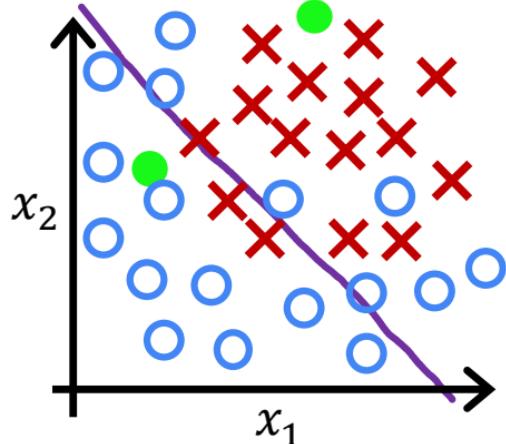


overfit

- Fits the training set extremely well

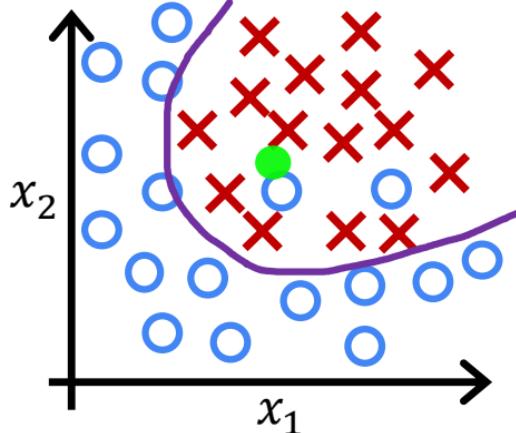
high variance

# Classification



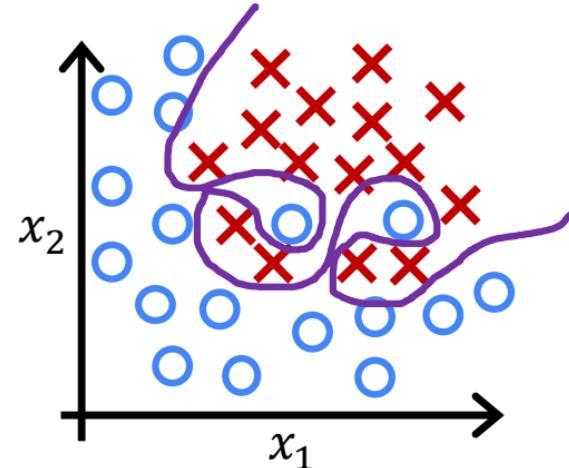
$$z = w_1 x_1 + w_2 x_2 + b$$
$$f_{\vec{w},b}(\vec{x}) = g(z)$$

$g$  is the sigmoid function  
underfit      high bias



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

just right

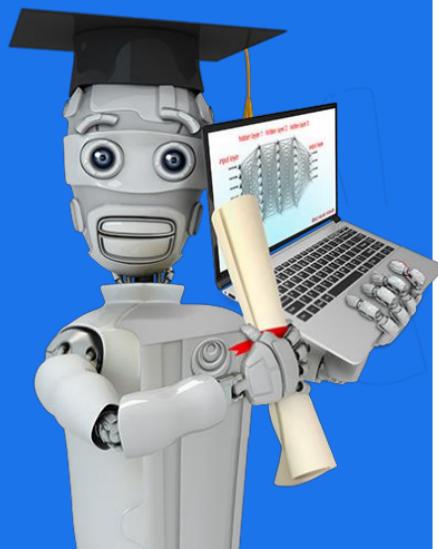


$$z = w_1 x_1 + w_2 x_2 + w_3 x_1^2 x_2 + w_4 x_1^2 x_2^2 + w_5 x_1^2 x_2^3 + w_6 x_1^3 x_2 + \dots + b$$

overfit

Stanford  
ONLINE

DeepLearning.AI

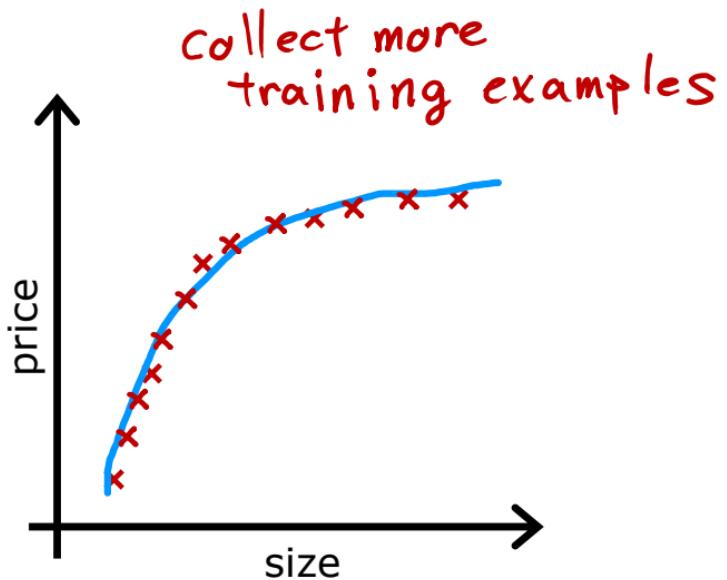
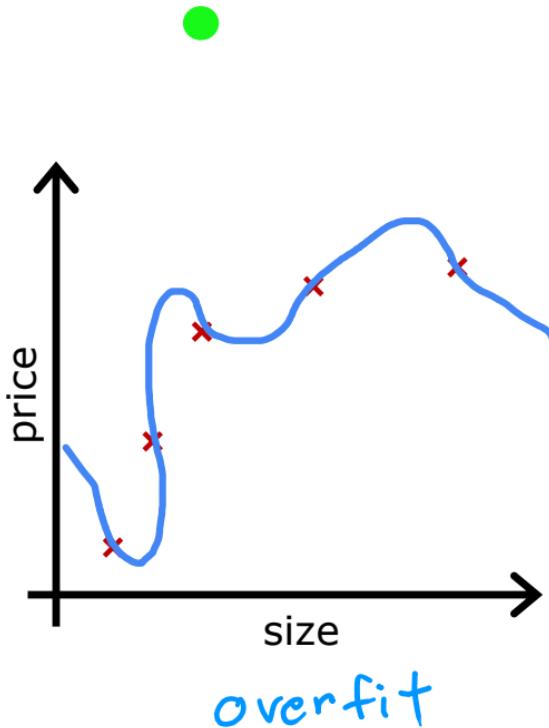


# Regularization to Reduce Overfitting

---

## Addressing Overfitting

# Collect more training examples



# Select features to include/exclude

size	bedrooms	floors	age	avg income	...	distance to coffee shop	price
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$		$x_{100}$	$y$

all features



insufficient data

↓  
overfit

selected features

size  
bedrooms  
age  
just right  
feature selection

course 2

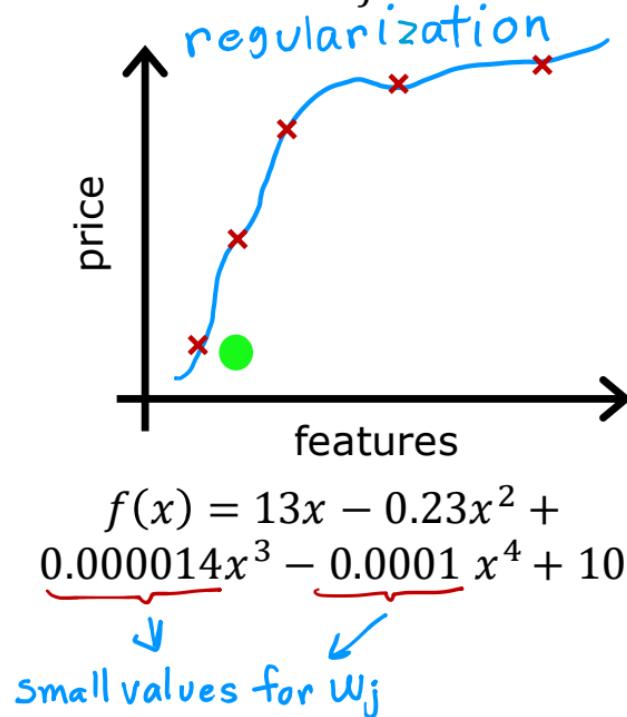
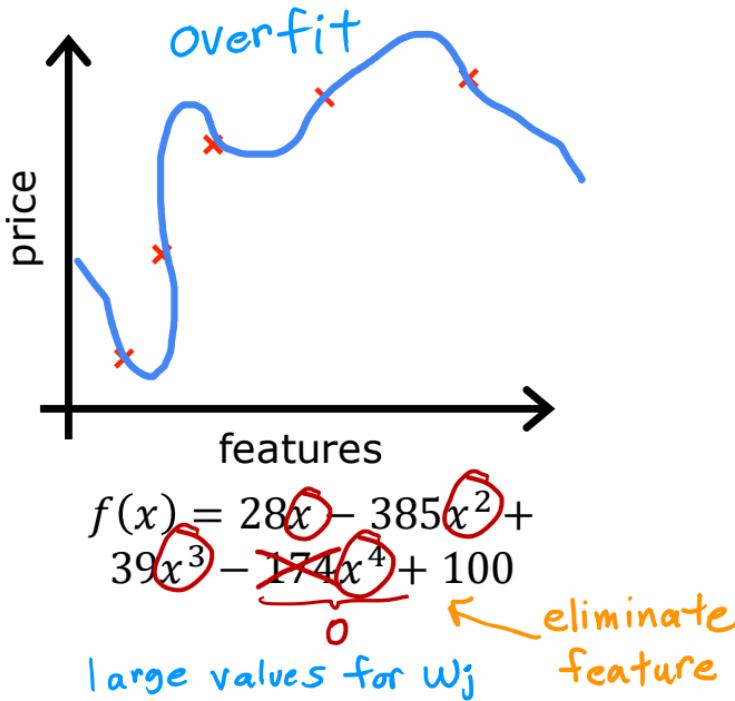
disadvantage



useful features could be lost

# Regularization

Reduce the size of parameters  $w_j$



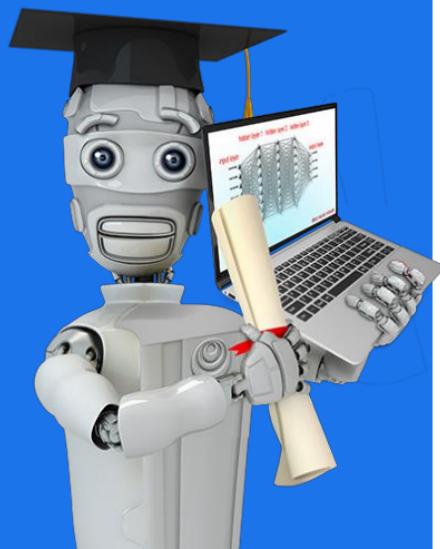
# Addressing overfitting

## Options

1. Collect more data
2. Select features
  - Feature selection *in course 2*
3. Reduce size of parameters
  - “Regularization” *next videos!*
    - Andrew Ng uses regularization all the time.

Stanford  
ONLINE

DeepLearning.AI

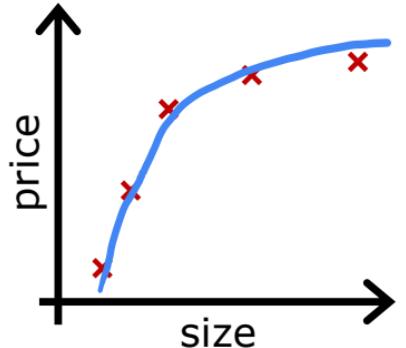


# Regularization to Reduce Overfitting

---

## Cost Function with Regularization

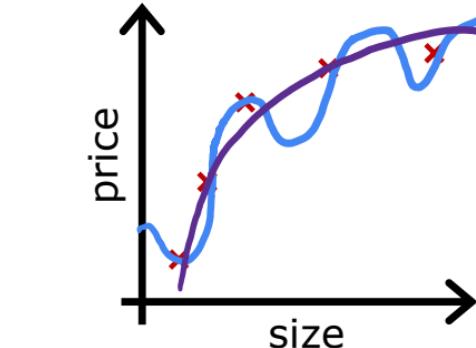
# Intuition



$$w_1x + w_2x^2 + b$$

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

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



$$w_1x + w_2x^2 + \cancel{w_3x^3} + \cancel{w_4x^4} + b$$

$\approx 0$        $\approx 0$

$$+ 1000 \underbrace{w_3^2}_{0.001} + 1000 \underbrace{w_4^2}_{0.002}$$

It penalizes the model  
for  $w_3$  and  $w_4$   
being large.

# Regularization

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

simpler model

$$w_3 \approx 0$$

less likely to overfit

$$w_4 \approx 0$$

size	bedrooms	floors	age	avg income	...	distance to coffee shop	price
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$		$x_{100}$	$y$

$w_1, w_2, \dots, w_{100}, b$

$n$  features

$n = 100$

Penalize all  $n$  features  $w_1, \dots, w_n$ .

regularization term

Andrew doesn't regularize  $b$  though.

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

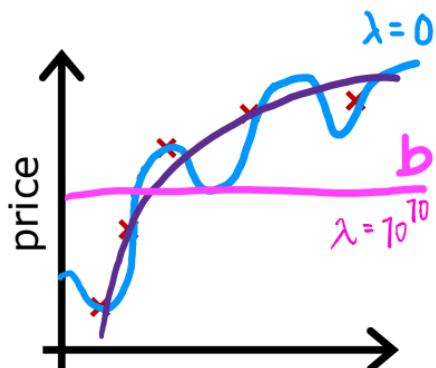
$\lambda > 0$

# Regularization

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

fit data      ↗ Keep  $w_j$  small

$\lambda$  balances both goals



choose  $\lambda = 10^{10}$

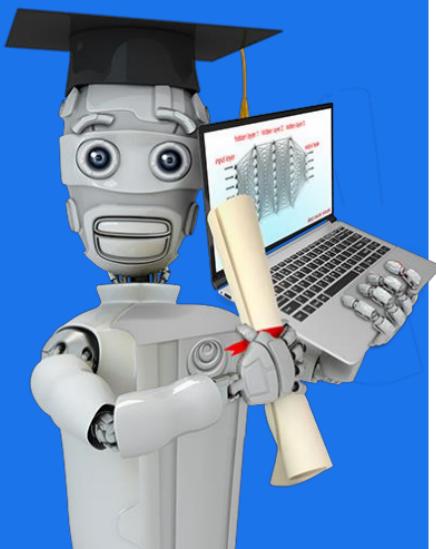
$$f_{\vec{w}, b}(\vec{x}) = \underbrace{w_1 x}_\approx + \underbrace{w_2 x^2}_\approx + \underbrace{w_3 x^3}_\approx + \underbrace{w_4 x^4}_\approx + b$$

$$f(x) = b$$

choose  $\lambda$  wisely

Stanford  
ONLINE

DeepLearning.AI



# Regularization to Reduce Overfitting

---

## 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

# 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  $j = 1, \dots, n$

# 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  $j = 1 \dots n$

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

usual update

$\star$  shrink  $w_j$

$$\alpha \frac{\lambda}{m}$$
$$0.01 \frac{1}{50} = 0.0002$$

$$w_j \underbrace{(1 - 0.0002)}_{0.9998}$$

$\star$

# How we get the derivative term (optional)

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

Stanford  
ONLINE

DeepLearning.AI

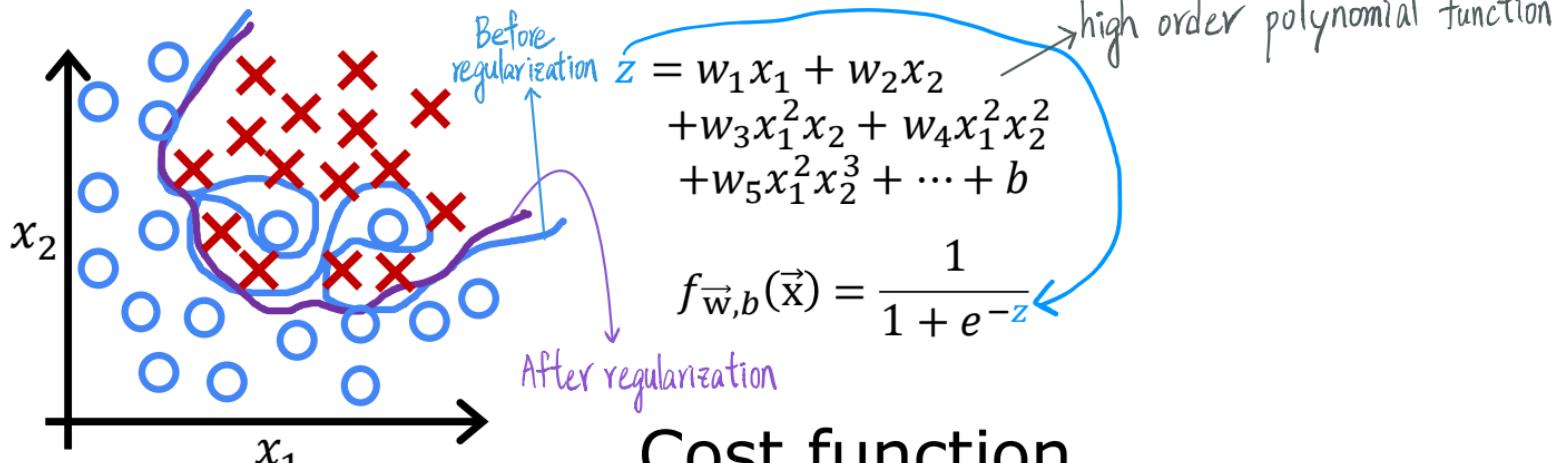


# Regularization to Reduce Overfitting

---

## Regularized Logistic Regression

# Regularized logistic regression



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

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

# 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 {

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

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

}

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  $b$