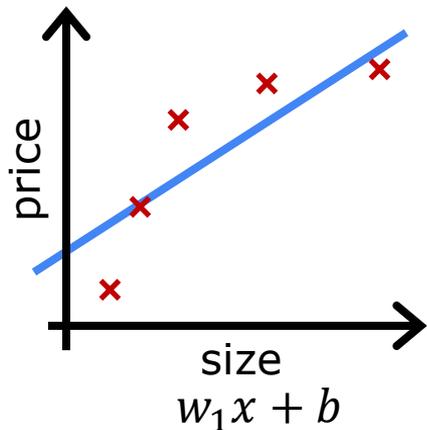# Regularization to Reduce Overfitting

## The Problem of Overfitting
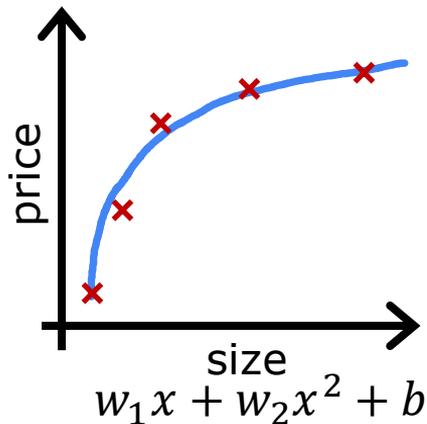
# Regression example

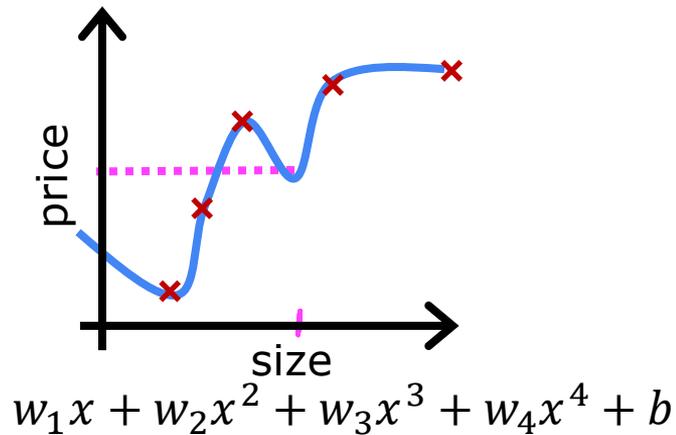

$$w_1 x + b$$

Underfit

- Does not fit the training set well

high bias

$$w_1 x + w_2 x^2 + b$$

- Fits training set pretty well

generalization

$$w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

overfit

- Fits the training set extremely well

high variance

# Classification



$$z = w_1 x_1 + w_2 x_2 + b$$
$$f_{\overrightarrow{\mathrm{w}},b}(\vec{\mathrm{x}}) = g(z)$$
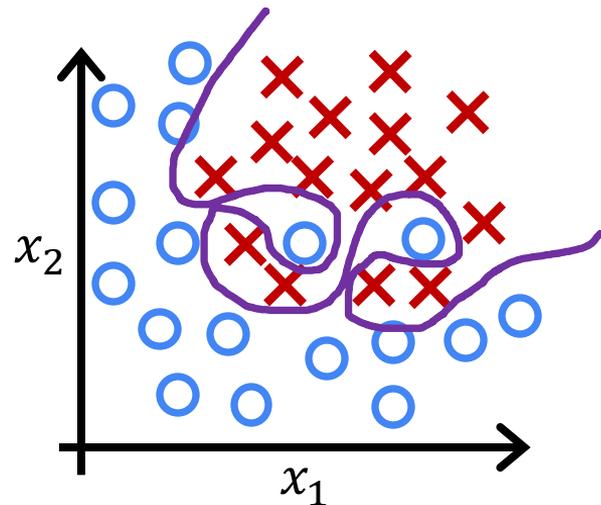
$g$ is the sigmoid function

underfit    high bias

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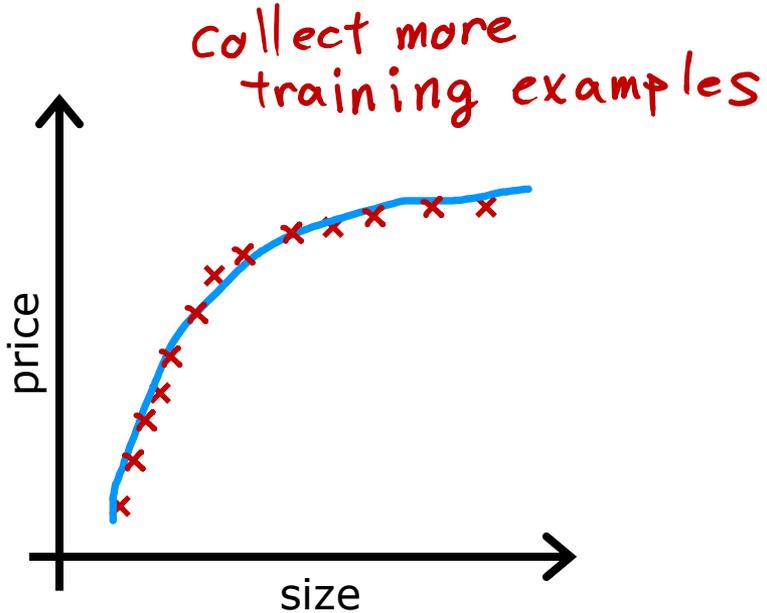$$+w_5 x_1 x_2 + b$$

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

Overfit

# 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

disadvantage

↓

useful features could be lost

# Regularization

Reduce the size of parameters $w_j$



Overfit

$$f(x) = 28x - 385x^2 + 39x^3 - \cancel{174x^4} + 10$$

0

large values for $w_j$

eliminate feature

regularization

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

small values for $w_j$

# Addressing overfitting

Options

1. Collect more data

2. Select features
   – Feature selection

3. Reduce size of parameters
   – "Regularization"

# Regularization to Reduce Overfitting

## Cost Function with Regularization

# Intuition



$$w_1 x + w_2 x^2 + b$$

$$w_1 x + w_2 x^2 + \underbrace{w_3 x^3}_{\approx 0} + \underbrace{w_4 x^4}_{\approx 0} + b$$

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

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

# Regularization

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

simpler model

less likely to overfit

$w_3 \hat{\approx} 0$

$w_4 \approx 0$

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

n features

n = 100

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

regularization term

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

"lambda"

regularization parameter

# Regularization

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

fit data

$\lambda$ balances both goals

← Keep $w_j$ small

choose $\lambda = 10^{10}$

$$f_{\vec{w},b}(\vec{x}) = \underset{\approx 0}{w_1 x} + \underset{\approx 0}{w_2 x^2} + \underset{\approx 0}{w_3 x^3} + \underset{\approx 0}{w_4 x^4} + b$$

$f(x) = b$

Choose $\lambda$



price

$\lambda = 0$

b

# Regularization to Reduce Overfitting

## Regularized Linear Regression

# Regularized linear regression

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

## 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

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

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

# Implementing gradient descent

repeat {

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

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

} simultaneous update    $j = 1...n$

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

$\underbrace{\phantom{w_j}}_{\text{shrink } w_j}$

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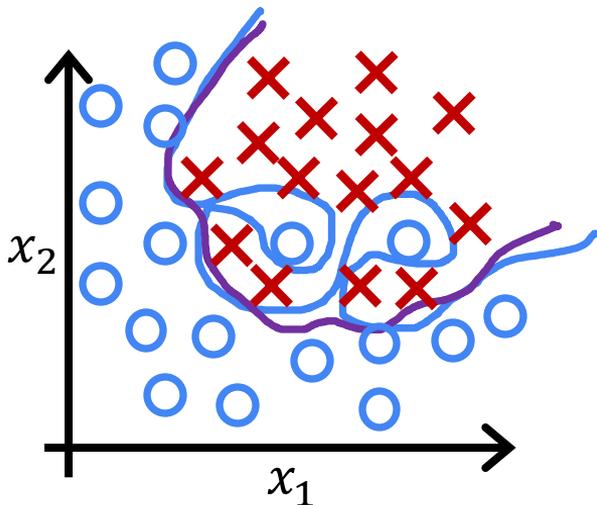
$0.01 \dfrac{1}{50} = 0.0002$

$w_j (1 - 0.0002)$

$0.9998$

# How we get the derivative term (optional)

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

$$\underbrace{\vec{w} \cdot \vec{x}^{(i)} + b}$$

$$= \frac{1}{2m} \sum_{i=1}^{m} \left[ \left( \vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)} \right) 2\, x_j^{(i)} \right] + \frac{\lambda}{2m} 2\, w_j \qquad \text{No } \sum_{j=1}^{n}$$

$$= \frac{1}{m} \sum_{i=1}^{m} \left( \underbrace{\left( \vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)} \right)}_{f(\vec{x})} x_j^{(i)} \right] + \frac{\lambda}{m} w_j$$

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

# Regularization to Reduce Overfitting

Regularized Logistic Regression

# Regularized logistic regression



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

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

## Cost function

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

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

# Regularized logistic regression

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

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

## Gradient descent

repeat {

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

$j = 1 \ldots n$

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

}

Looks same as for linear regression!

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

logistic regression

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