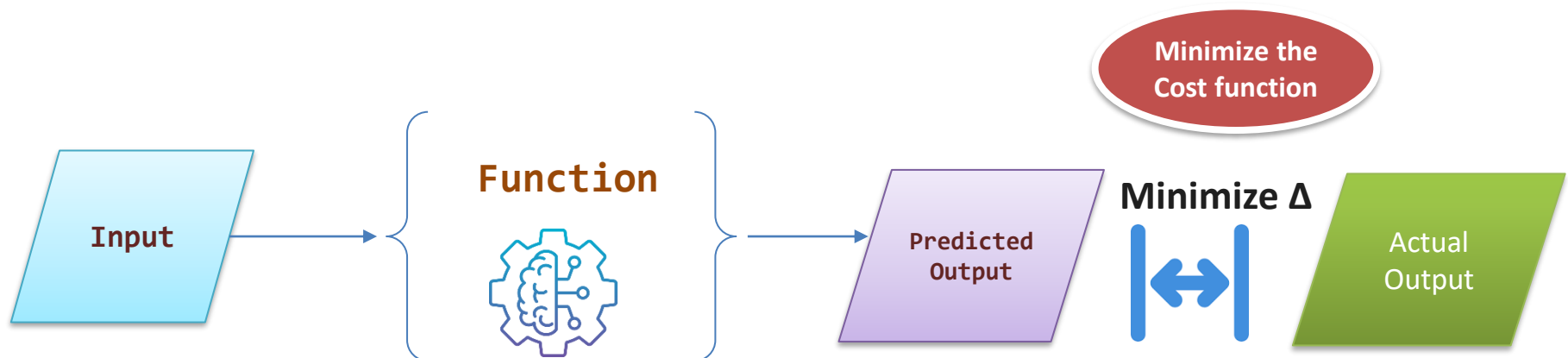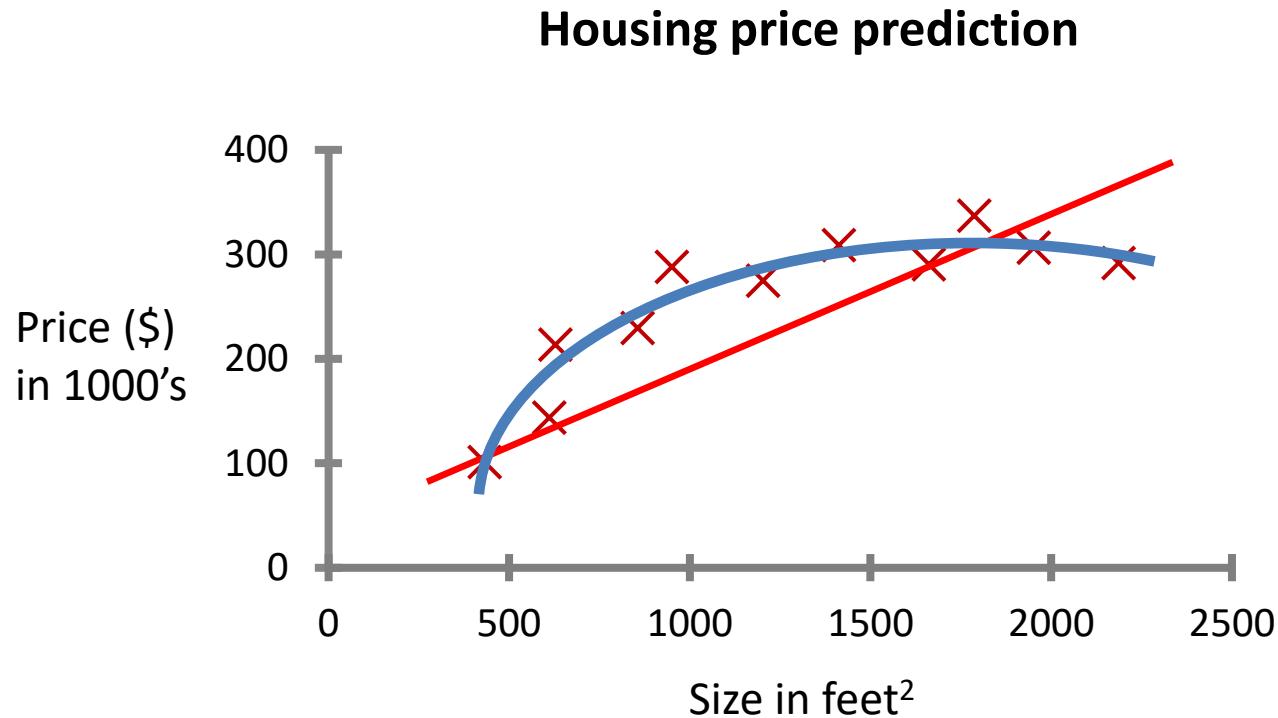# Regression

# ML: learn a **Function** that minimizes the cost

- Start with random function parameters
- Repeat intelligent guessing/approximation of the Function parameters such that the difference between the Predicted Output the Actual Output is reduced
  - i.e., minimize a Cost function a.k.a loss, or error function

# Linear Regression with One Variable

**Housing price prediction**



**Regression**: Predict continuous output value (price)

**Training set of housing prices**

| Size in feet² ($x$) | Price ($) in 1000's ($y$) |
|:---:|:---:|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

Notation:

$m$ = Number of training examples

$x$ = "input" variable / features

$y$ = "output" variable / "target" variable
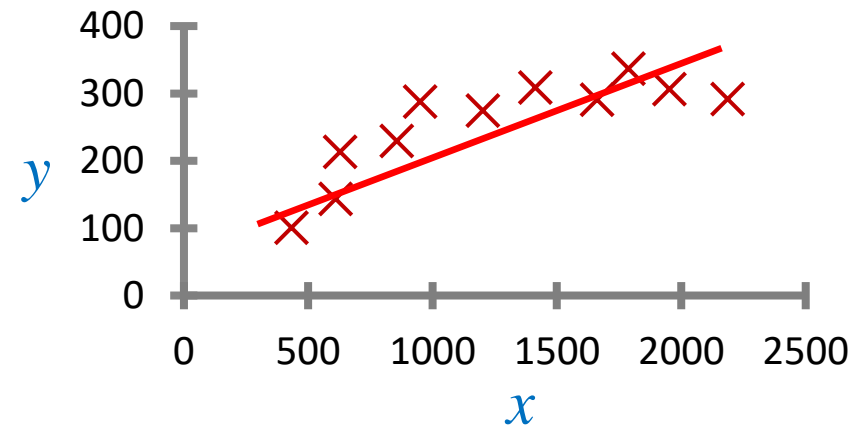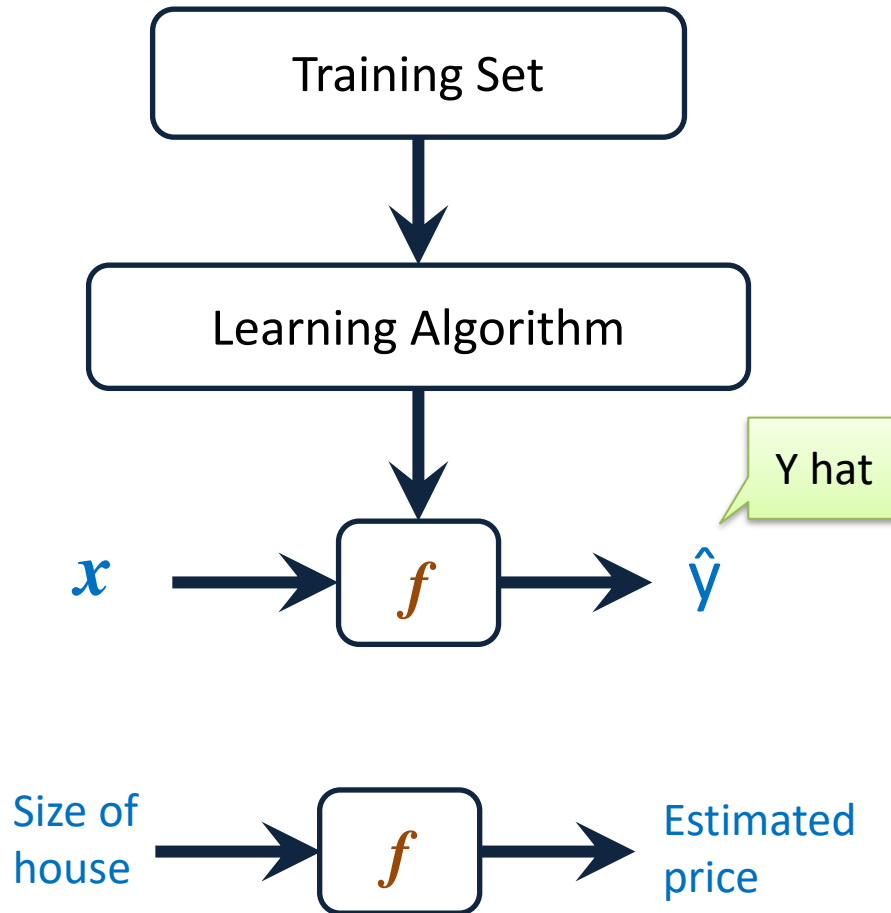
$x^{(1)} = 2104$

$x^{(2)} = 1416$

$y^{(1)} = 460$

$(x^{(i)}, y^{(i)})$ – the $i^{th}$ training example

# How do we represent $f$ ?

Training Set

↓

Learning Algorithm

↓

$x$ → $f$ → $\hat{y}$

**Y hat**

Size of house → $f$ → Estimated price

$$f(x) = wx + b$$

$w, b$ are parameters (coefficients) to learn from the training set



Linear regression with one variable
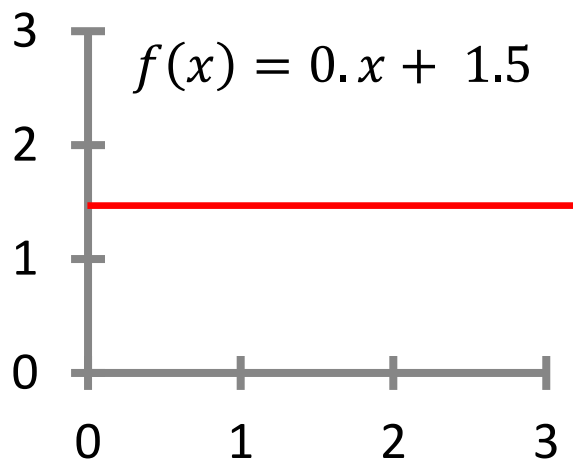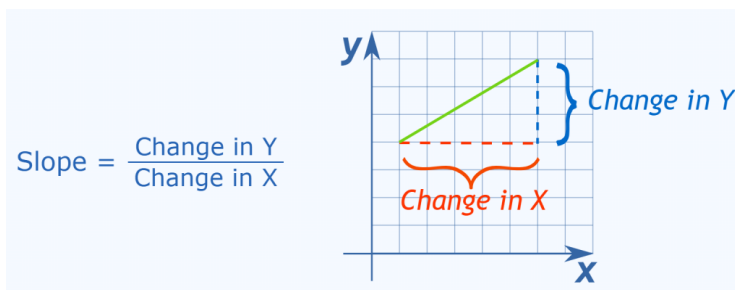**Univariate linear regression**

Given a training set, **learn a function** $f$ so that $f(x)$ is a "good" predictor for the corresponding value of y (i.e. minimize the error between predicted and actual values)

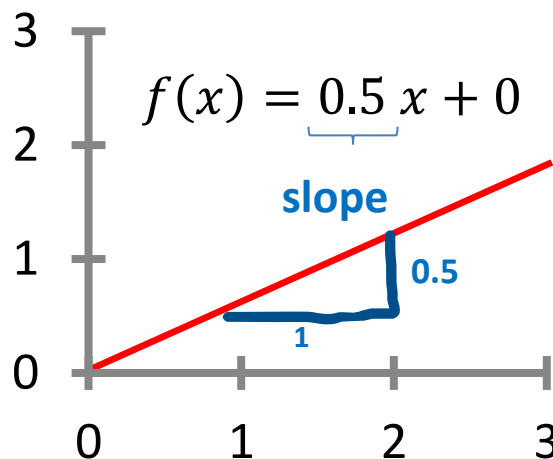# Univariate Linear Regression - Model Representation

$$f(x) = wx + b$$

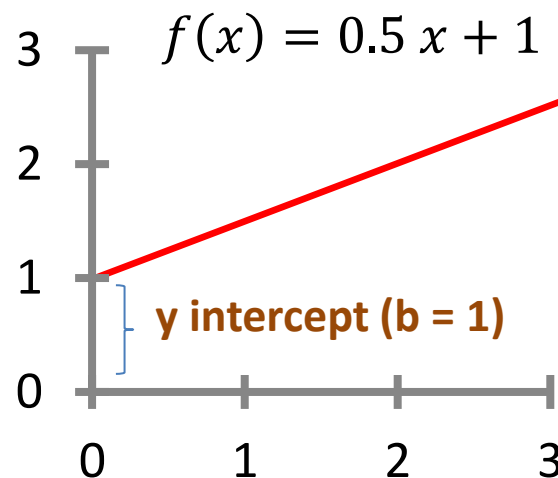- $w$ is the slope of the line
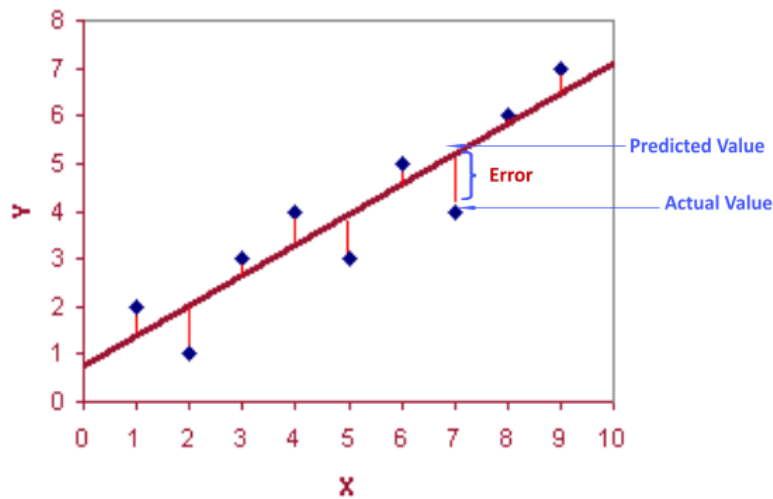- $b$ is the y-intercept of the line

## How to choose $w$ and $b$ ?



$$\text{Slope} = \frac{\text{Change in Y}}{\text{Change in X}}$$

$f(x) = 0.\,x + \,1.5$

$f(x) = 0.5\,x + 0$

slope

0.5

1

$f(x) = 0.5\,x + 1$

y intercept (b = 1)

$w = 0$
$b = 1.5$

$w = 0.5$
$b = 0$

$w = 0.5$
$b = 1$

**Idea**: Choose $w$ and $b$ so that $f(x)$ is close to $y$ for our training examples $(x, y)$

Find $w, b$:
$\hat{y}^{(i)}$ is close to $y^{(i)}$ for all $(x^{(i)}, y^{(i)})$

**Cost (mean squared error)**
Function:
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^{m} (f(x^{(i)}) - y^{(i)})^2$$

Goal: $\underset{w, b}{\text{minimize}} \, J(w, b)$

With **m** = number of training examples

Function:

$$f(x) = wx + b$$

Parameters:

$$w, b$$

Cost Function:

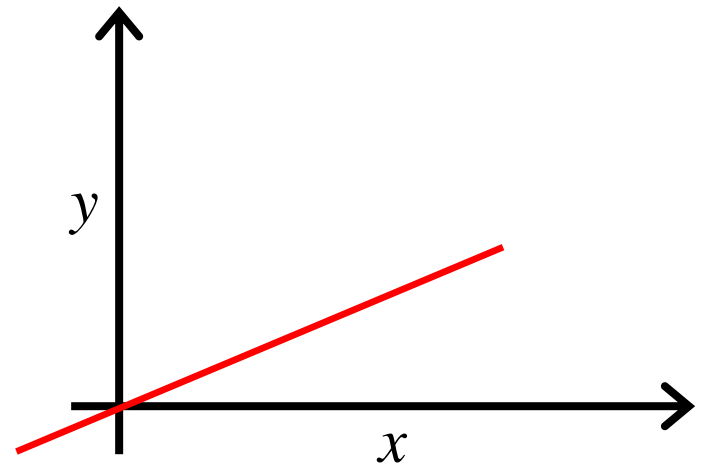$$J(w, b) = \frac{1}{2m} \sum_{i=1}^{m} (f(x^{(i)}) - y^{(i)})^2$$

Goal:   $\underset{w, b}{\text{minimize}} \ J(w, b)$
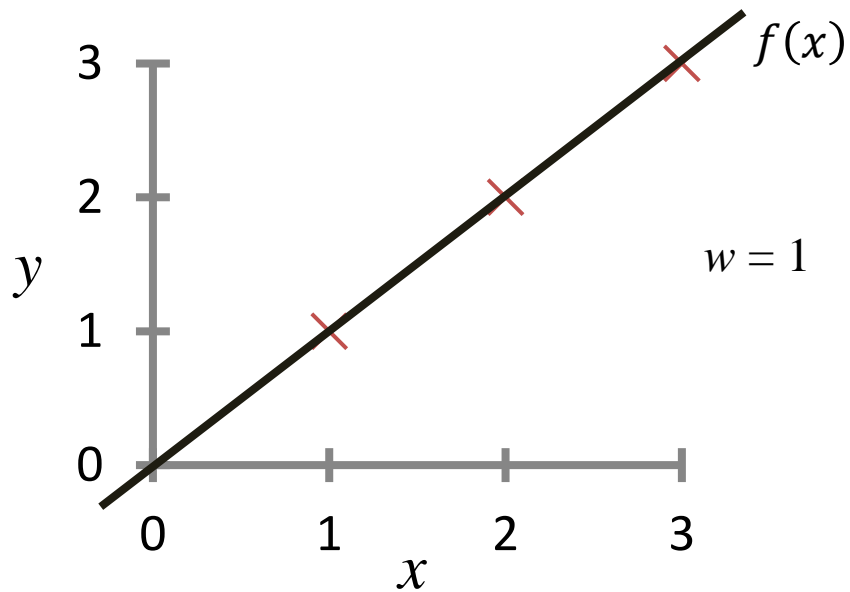
$$f(x) = wx$$

$$w$$

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

$\underset{w}{\text{minimize}} \ J(w)$

# $f(x) = wx$

$J(w)$

(function of the parameter $w$)



$f(x)$

$w = 1$

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

$$= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$$

# $f(x) = wx$

$J(w)$

$f(x)$

$w = 0.5$

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

$$= \frac{1}{2m}((0.5 - 1)^2 + (1\text{-}2)^2 + (1.5\text{-}3)^2)$$

$$= \frac{1}{2\times3}(3.5) = \frac{3.5}{6} = 0.58$$

$$f(x) = wx$$



$$w = 0$$

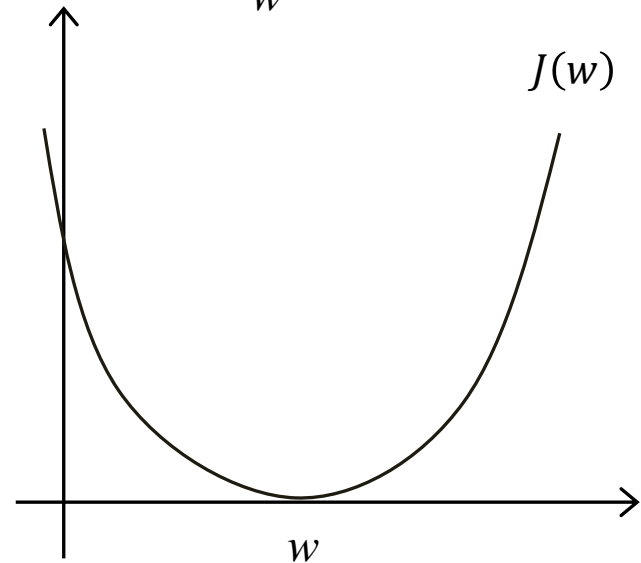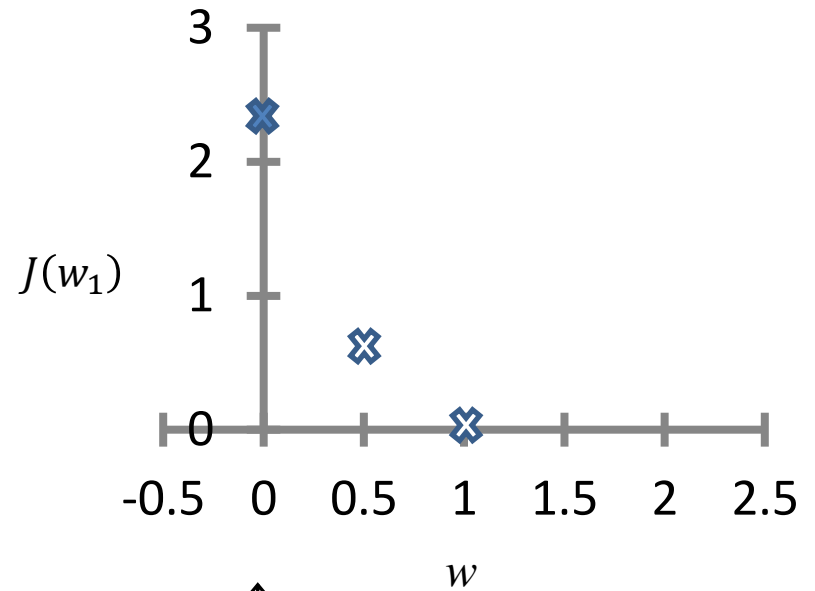$$J(w) = \frac{1}{2m} \sum_{i=1}^{m} (f(x^{(i)}) - y^{(i)})^2$$

$$= \frac{1}{2m} (1^2 + 2^2 + 3^2)$$

$$= \frac{1}{2 \times 3} (14) = \frac{14}{6} = 2.3$$

$$J(w)$$



$J(w_1)$

$$J(w)$$

$$f(x) = wx + b$$

(for fixed $w, b$ this is a function of $x$)

$$J(w, b)$$

(function of the parameters $w, b$)



Price ($) in 1000's — Size in feet$^2$ ($x$)

$$f(x) = 50 + 0.06x$$

$$f(x) = wx + b$$

(for fixed $w$, $b$ this is a function of $x$)

$$J(w, b)$$

(function of the parameters $w$, $b$)

Contour Plot

$$f(x) = wx + b$$

(for fixed $w$, $b$ this is a function of $x$)

$$J(w, b)$$

(function of the parameters $w$, $b$)

$$f(x) = wx + b$$

(for fixed $w$, $b$ this is a function of $x$)

$$J(w, b)$$

(function of the parameters $w$, $b$)

$$f(x) = wx + b$$
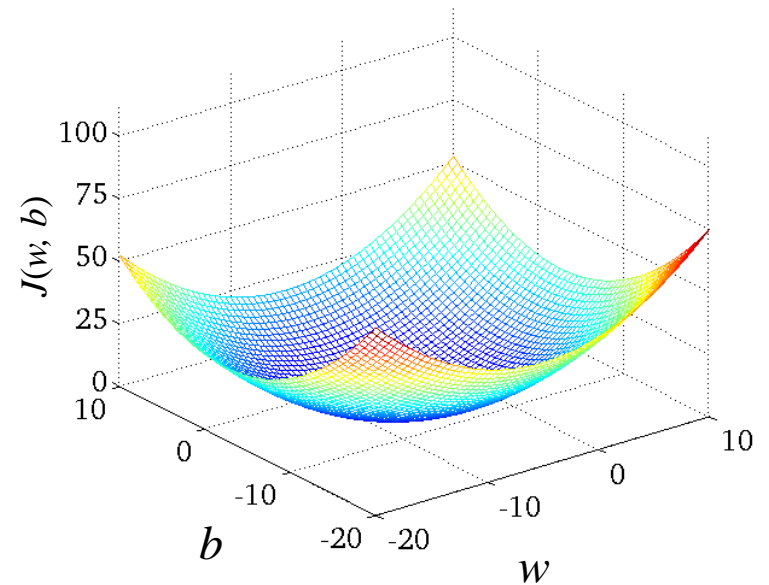
(for fixed $w$, $b$ this is a function of $x$)

$$J(w, b)$$

(function of the parameters $w$, $b$)

$$f(x) = wx + b$$

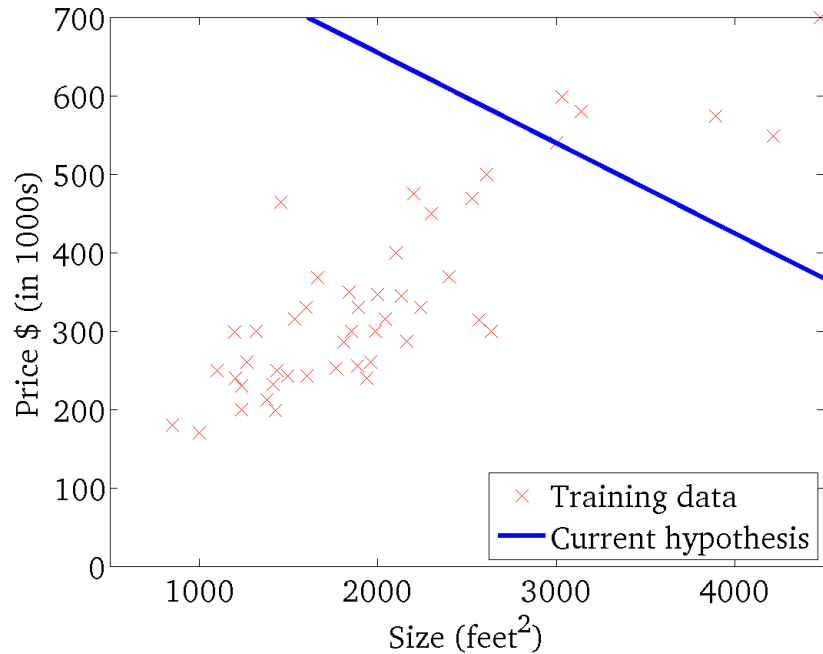(for fixed $w$, $b$ this is a function of $x$)

$$J(w, b)$$

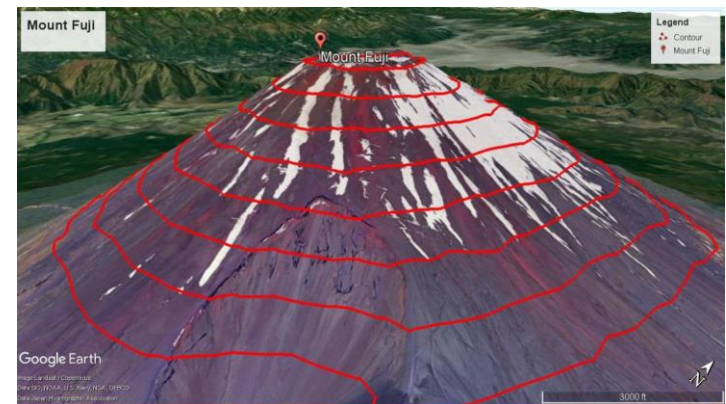(function of the parameters $w$, $b$)

$$f(x) = wx + b$$

(for fixed $w, b$ this is a function of $x$)

$$J(w, b)$$

(function of the parameters $w, b$)

$$f(x) = wx + b$$

(for fixed $w$, $b$ this is a function of $x$)

$$J(w, b)$$

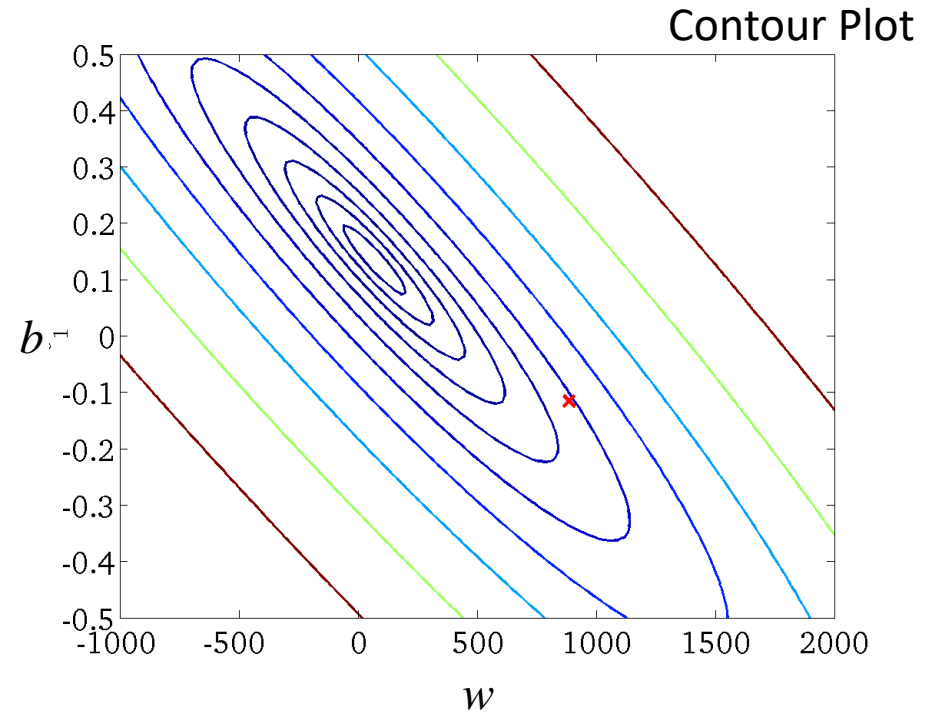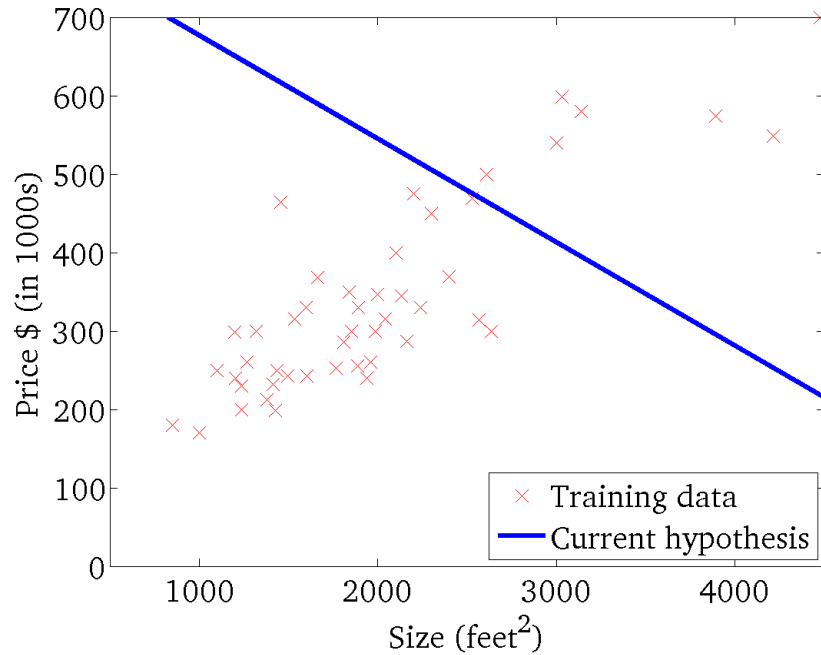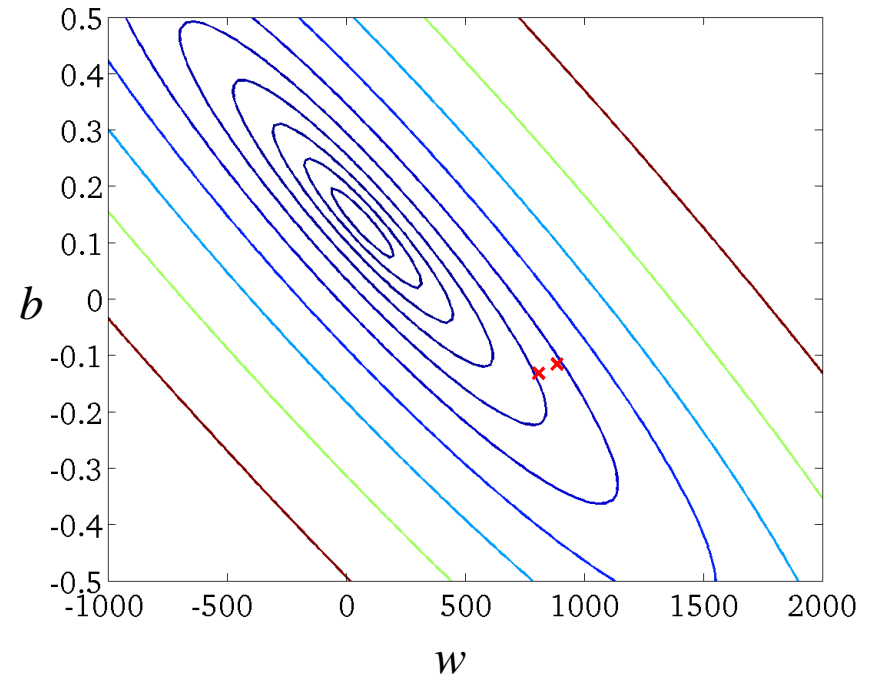(function of the parameters $w$, $b$)

$$f(x) = wx + b$$

(for fixed $w, b$ this is a function of $x$)
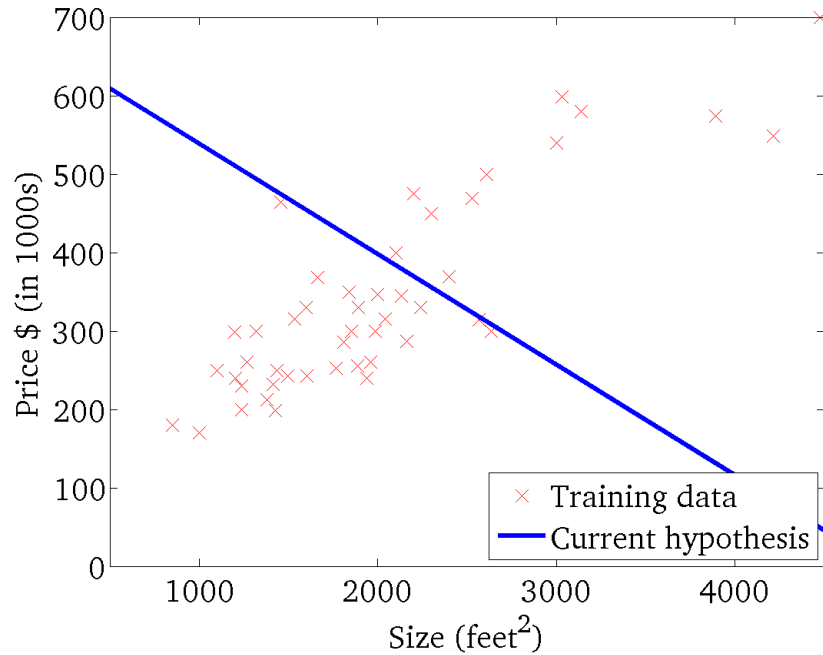
$$J(w, b)$$

(function of the parameters $w, b$)

$$f(x) = wx + b$$

(for fixed $w$, $b$ this is a function of $x$)

$$J(w, b)$$

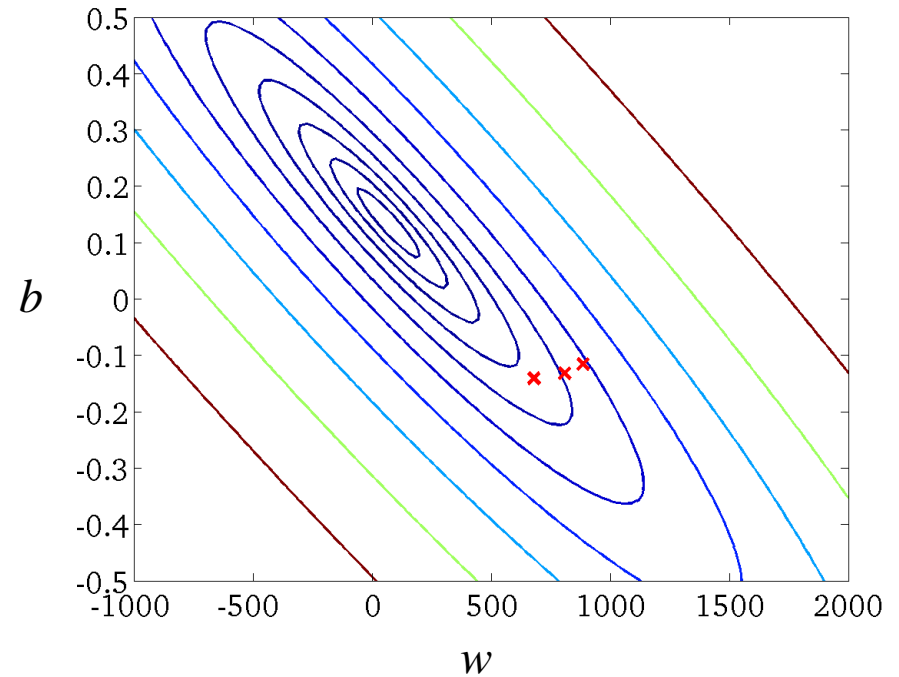(function of the parameters $w$, $b$)

# Linear Regression - Gradient decent

- Gradient decent

Have a cost function $J(w, b)$

Want to find w and be $\min\limits_{w,b} J(w, b)$

**Outline:**

- Start with some $w, b$ (say 0, 0)

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

until we hopefully end up at a minimum

$J(w, b)$

$w$

$b$

$J(w, b)$

$w$

$b$

# Gradient descent algorithm

1. Initialize the values of $w$ and $b$ to some arbitrary values

2. Calculate the predicted values of $y$ using the current values of $w$ and $b$

3. Calculate the gradients of the cost function with respect to $w$ and $b$

4. Update the values of $w$ and $b$ using the gradients and a learning rate

5. Repeat steps 2-4 until convergence (i.e., until the cost function converges to a minimum)

# Gradient descent algorithm

Repeat until convergence

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

Learning rate

Derivative

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

Simultaneously update w and b

Correct: Simultaneous update

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

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

$$w = tmp\_w$$

$$b = tmp\_b$$

Incorrect

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

$$w = tmp\_w$$

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

$$b = tmp\_b$$

# Derivative 101

- Source

- **Derivatives**: it is all about slope!

$$\text{Slope} = \frac{\text{Change in Y}}{\text{Change in X}}$$

We can find
an **average** slope between
two points

$$\text{average slope} = \frac{24}{15}$$

# But how do we find the slope at a point?

- There is nothing to measure!
  slope 0/0 = ????
- But with derivatives we use a small difference …

… then have it shrink towards zero

- Fill in this slope formula: $\frac{\Delta y}{\Delta x} = \frac{f(x+\Delta x) - f(x)}{\Delta x}$
- Simplify it as best we can
- Then make $\Delta x$ shrink towards zero.

$$\text{slope} = \frac{0}{0} = ???$$

# Derivative Example - $f(x) = x^2$

The slope formula is: $\dfrac{f(x+\Delta x) - f(x)}{\Delta x}$

Use $f(x) = x^2$: $\dfrac{(x+\Delta x)^2 - x^2}{\Delta x}$

Expand $(x+\Delta x)^2$ to $x^2+2x\,\Delta x+(\Delta x)^2$: $\dfrac{x^2 + 2x\,\Delta x + (\Delta x)^2 - x^2}{\Delta x}$

Simplify ($x^2$ and $-x^2$ cancel): $\dfrac{2x\,\Delta x + (\Delta x)^2}{\Delta x}$

Simplify more (divide through by $\Delta x$): $2x + \Delta x$

Then, **as Δx heads towards 0** we get: $2x$

Result: the derivative of $x^2$ is **2x**    $\dfrac{d}{dx}x^2 = 2x$

- In other words, the slope at x is 2x

# Interpretation of Derivative

- So what does $\frac{d}{dx}x^2 = 2x$ mean?

- It means that, for the function $x^2$, the slope or "rate of change" at any point is **2x**

- So when **x=2** the slope is **2x = 4**

- Or when **x=5** the slope is **2x = 10**, and so on

# Gradient descent algorithm

repeat until convergence {

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

}



*unsplash.com/photos/3m6vbzY69s4*

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

$J(w)$

Positive slope/derivative

$w$

$J(w)$

Negative slope/derivative

$w$

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

If α is too small, gradient descent can be slow



If α is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge

Gradient descent can converge to a local minimum, even with the learning rate α fixed

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

As we approach a local minimum, gradient descent will automatically take smaller steps (the slope gets smaller). So, no need to decrease α over time

# Linear Regression with One Variable

## 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 {

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

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

}

# "Batch" Gradient Descent

"Batch": Each step of gradient descent uses all the training examples, $m$

| | $x$ size in feet² | $y$ price in $1000's |
|---|---|---|
| (1) | 2104 | 400 |
| (2) | 1416 | 232 |
| (3) | 1534 | 315 |
| (4) | 852 | 178 |
| ... | ... | ... |
| (47) | 3210 | 870 |

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

# Linear Regression with multiple variables

**Multiple features (variables)**

| Size (feet²) | Number of bedrooms | Number of floors | Age of home (years) | Price ($1000) |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| … | … | … | … | … |

**Notation:**

$n$ = number of features

$x_j$ = $j^{th}$ feature

$\vec{x}^{(i)}$ = features of $i^{th}$ training example

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

Hypothesis:

Previously:   $f(x) = wx + b$

 Now:   Multivariate linear regression.

$$f(x) = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$$

**Example**   $f(x) = 0.1 x_1 + 4 x_2 + 10 x_3 + -2 x_4 + 80$

size   #bedrooms   #floors   years   base price

Parameters:   $w_0, w_1, \ldots, w_n$

Cost function:
$$J(w_0, w_1, \ldots, w_n) = \frac{1}{2m} \sum_{i=1}^{m} (f(x^{(i)}) - y^{(i)})^2$$

Gradient descent:

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_0, \ldots, w_n)$$

}                    (simultaneously update for every $j = 0, \ldots, n$)

## Gradient Descent

Previously ($n=1$):

Repeat {
$$w = w - \alpha \frac{1}{m} \underbrace{\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial w}J(w)}$$

$$b = b - \alpha \frac{1}{m} \underbrace{\sum_{i=1}^{m}\left(f_w(x^{(i)}) - y^{(i)}\right) \cdot x^{(i)}}_{\frac{\partial}{\partial b}J(w)}$$
}

(simultaneously update $w_0$ and $w_1$)

New algorithm  $(n \geq 1)$:

Repeat {
$$w_j = w_j - \alpha \frac{1}{m} \underbrace{\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\frac{\partial}{\partial w_j}J(w)}$$

(simultaneously update $w_j$ for $j = 0, \dots, n$)
}

$$w_0 = w_0 - \alpha \frac{1}{m} \sum_{i=1}^{m}\left(f_w(x^{(i)}) - y^{(i)}\right) x_0^{(i)}$$

$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^{m}\left(f_w(x^{(i)}) - y^{(i)}\right) x_1^{(i)}$$

$$w_2 = w_2 - \alpha \frac{1}{m} \sum_{i=1}^{m}\left(f_w(x^{(i)}) - y^{(i)}\right) x_2^{(i)}$$

….

Formula to update **b** remains the same

# Gradient descent in practice:

- Feature Scaling
- Regularization
- Regression Evaluation

**Feature Scaling:** divide the input values by the range (i.e. the maximum value minus the minimum value) of the input variable, resulting in a new range of just 1.

The idea: Make sure features are on a similar scale. So that the gradient descent converges faster.

E.g. $x_1$ = size (0-2000 feet$^2$)

$x_2$ = number of bedrooms (1-5)

$$x_1 = \frac{\text{size (feet}^2)}{2000}$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$



Rule-of-thumb: Get every feature into approximately a $-1 \leq x_i \leq 1$ range, $-0.5 \leq x_i \leq 0.5$, or other similar small ranges.

# Mean normalization

- Replace $x_i$ to make features have approximately zero mean (Do not apply to $x_0 = 1$):

$$x_i := \frac{x_i - \mu_i}{s_i}$$
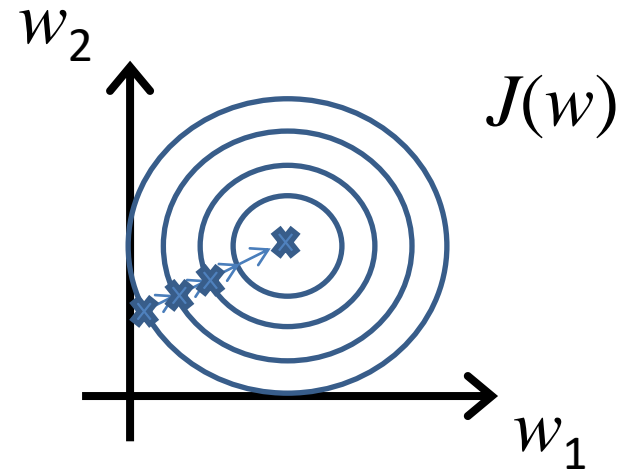
Where $\mu_i$ is the **average** of all the values for feature ($i$) (**in the training set**) and $s_i$ is the range of values ($max$ - $min$), or $s_i$ is the standard deviation.

- e.g.,

$x_1 = \frac{size - 1000}{2000}$  (average size of the houses is 1000, and ranges from 0 to 2000)

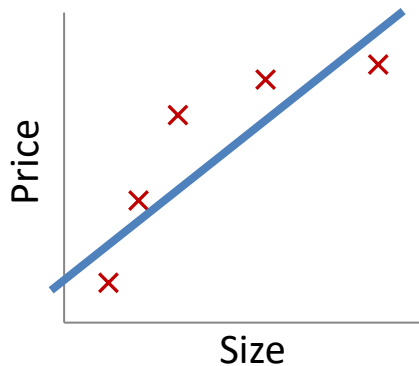$x_2 = \frac{\#bedrooms - 2}{4}$  (average # of bedrooms is 2, and the range is from 1 to 5)

$$-0.5 \le x_1 \ge 0.5, -0.5 \le x_2 \ge 0.5,$$

# Regularization

- The problem of overfitting

Example: Linear regression (housing prices)



$$w_0 + w_1 x$$

"underfit/high bias"

$$w_0 + w_1 x + w_2 x^2$$

"just right"

$$w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4$$

"overfit/high variance"

**Overfitting:** If we have too many features, the learned hypothesis may fit the training set very well ( $J(w) = \frac{1}{2m} \sum_{i=1}^{m} (h_w(x^{(i)}) - y^{(i)})^2 \approx 0$ ) but fail to generalize to new examples (predict prices on new examples).

## Addressing overfitting:

$x_1 =$ size of house
$x_2 =$ no. of bedrooms
$x_3 =$ no. of floors
$x_4 =$ age of house
$x_5 =$ average income in neighborhood
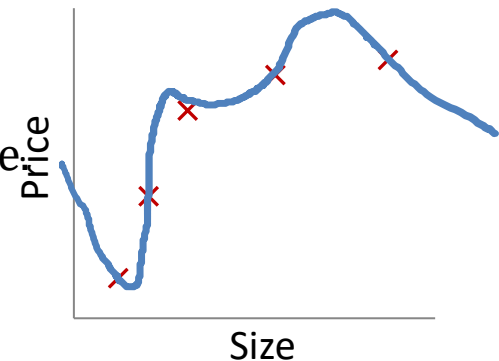$x_6 =$ kitchen size
$\vdots$

$x_{100}$



## Addressing overfitting:

Options:

1. Reduce number of features.
   — Manually select which features to keep.
   — Use feature selection algorithm.
2. Regularization.
   — Keep all the features, but reduce magnitude/values of parameters $w_j$
   — Works well when we have a lot of features, each of which contributes a bit to predicting $y$ .

# Regularization

- Cost function

**Intuition**



$$w_0 + w_1 x + w_2 x^2$$

$$w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4$$

Suppose we penalize and make $w_3, w_4$ really small

$$J(w) = \frac{1}{2m} \sum_{i=1}^{m} (f_w(x^{(i)}) - y^{(i)})^2 + 1000 \, w_3^2 + 1000 \, w_4^2$$

$$w_3 \approx 0, \quad w_4 \approx 0$$

# Regularization

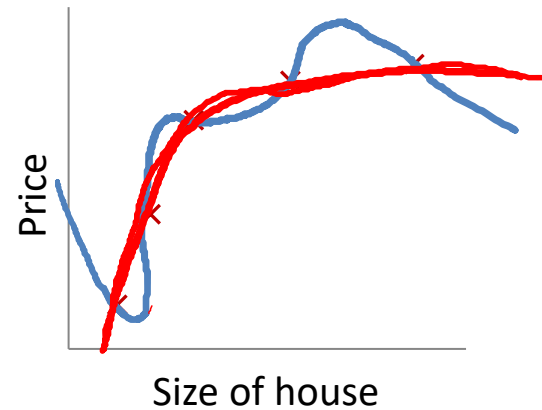Small values for parameters $w_0, w_1, \ldots, w_n$
- — "Simpler/smoother" hypothesis
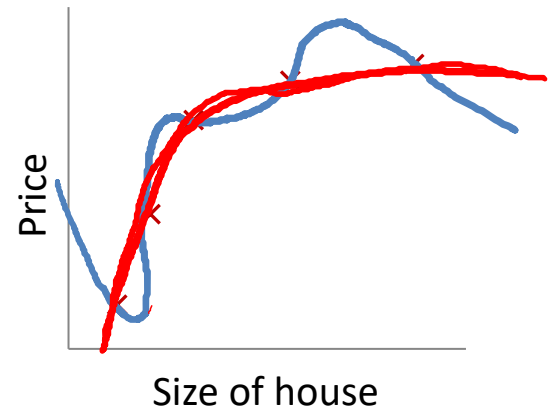- — Less prone to overfitting

Housing:
- — Features: $x_1, x_2, \ldots, x_{100}$
- — Parameters: $w_0, w_1, w_2, \ldots, w_{100}$

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

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

$$\min_{w} J(w)$$

Price

Size of house

In regularized linear regression, we choose $w$ to minimize

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

What if $\lambda$ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$ )?



Price

Size of house

"underfit"

$w_1 \approx 0, \quad w_2 \approx 0, \quad w_3 \approx 0, \quad w_4 \approx 0$

$w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4$

## Regularized linear regression

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

$$\min_{w} J(w)$$

## Gradient descent

$$\frac{\partial}{\partial w_0}J(w)$$

Repeat {

$$w_0 = w_0 - \alpha \frac{1}{m}\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})\, x_0^{(i)}$$

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

$$\frac{\partial}{\partial w_j}J(w) \text{ "Regularized"}$$

} $(j = \cancel{0}, 1, 2, 3, \ldots, n)$

$$w_j := w_j\left(1 - \alpha\frac{\lambda}{m}\right) - \alpha\frac{1}{m}\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})\, x_j^{(i)}$$

$$1 - \alpha\frac{\lambda}{m} < 1$$

# Regression Evaluation

- Performance measured by

  - Mean Squared Error (MSE)
    $$MSE = \frac{1}{n}\sum(y - \hat{y})^2$$

  - Root-Mean-Squared-Error (RMSE)
    $$RMSE = \sqrt{\frac{(y - \hat{y})^2}{n}}$$

  - Mean-Absolute-Error (MAE)
    $$MAE = \frac{1}{n}\sum|y - \hat{y}|$$

  - …others