# Linear Regression



https://xkcd.com/1007/

Slides from:

Andrew Ng

CSC411: Machine Learning and Data Mining, Winter 2018

Michael Guerzhoy and Lisa Zhang

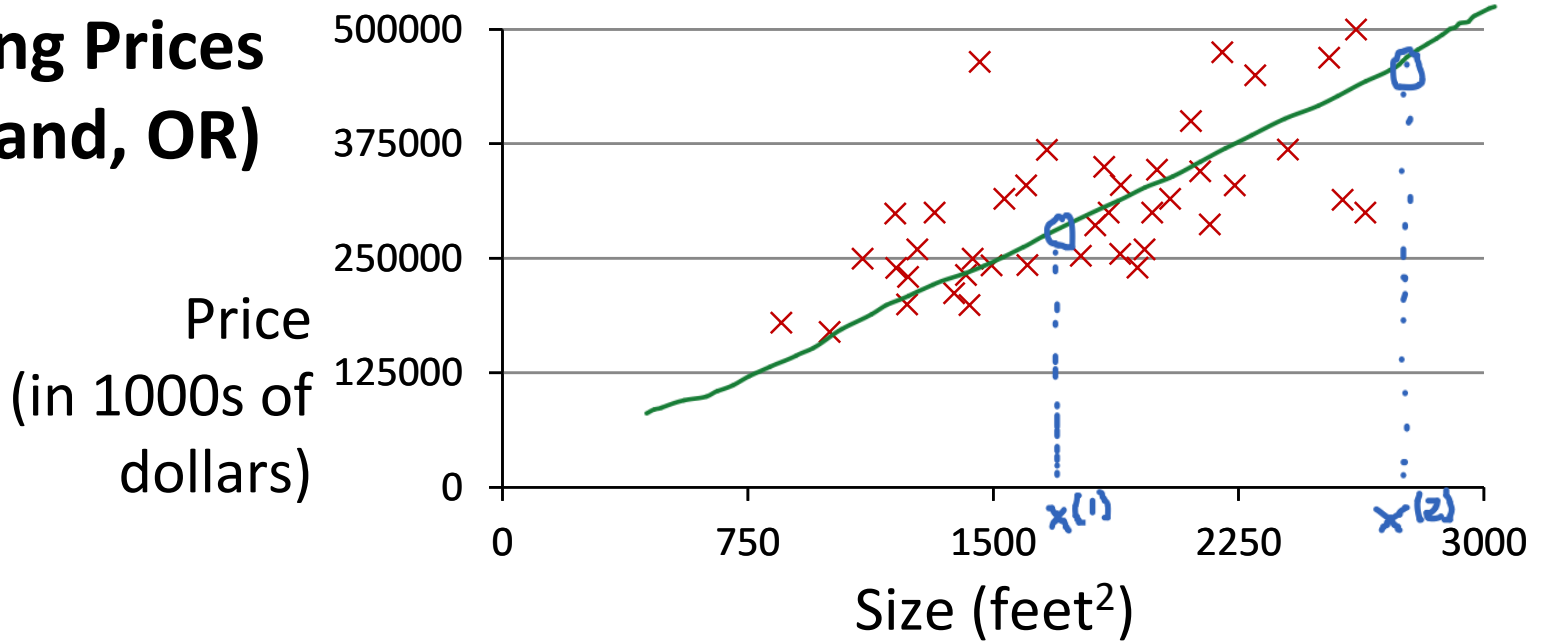| **Training set of housing prices (Portland, OR)** | **Size in feet² ($x$)** | **Price ($) in 1000's ($y$)** |
|---|---|---|
| | 2104 | 460 |
| | 1416 | 232 |
| | 1534 | 315 |
| | 852 | 178 |
| | … | … |

Notation:

**m** = Number of training examples
**x**'s = "input" variable / features
**y**'s = "output" variable / "target" variable
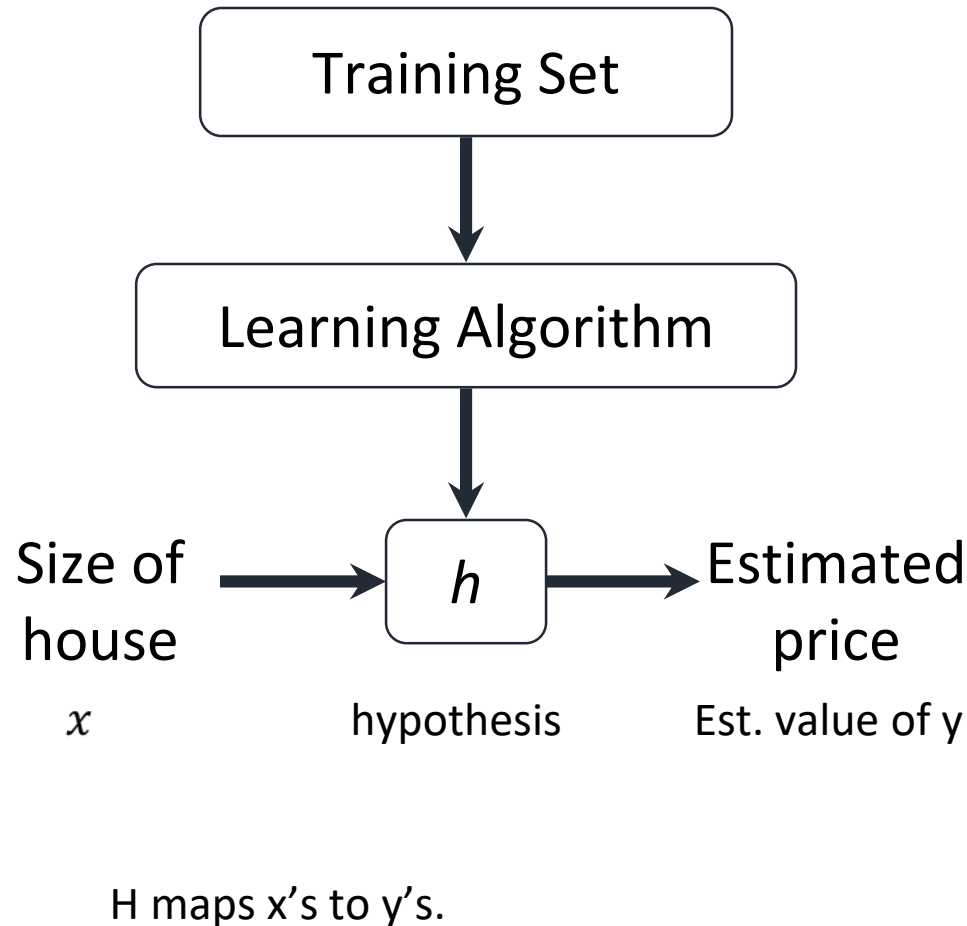
**Housing Prices (Portland, OR)**



Price (in 1000s of dollars)

Size (feet$^2$)

Supervised Learning

Given the "right answer" for each example in the data.

Regression Problem

Predict real-valued output

| Training Set |

$\downarrow$

| Learning Algorithm |

$\downarrow$

Size of house $\longrightarrow$ | $h$ | $\longrightarrow$ Estimated price

$x$           hypothesis          Est. value of y

H maps x's to y's.

## How do we represent $h$ ?

- We represent hypotheses about the data using the parameters $\theta = (\theta_0, \theta_1)$
- If the data is correctly predicted according to hypothesis $h_\theta$, then $y \approx h_\theta(x) = \theta_0 + \theta_1 x$
- The learning algorithm finds the best hypothesis $h_\theta$ for the training set
- We can then estimate the values of y for the test set using that $h_\theta$
- If $h_\theta(x)$ is a linear function of a real number x, this procedure is called linear regression.

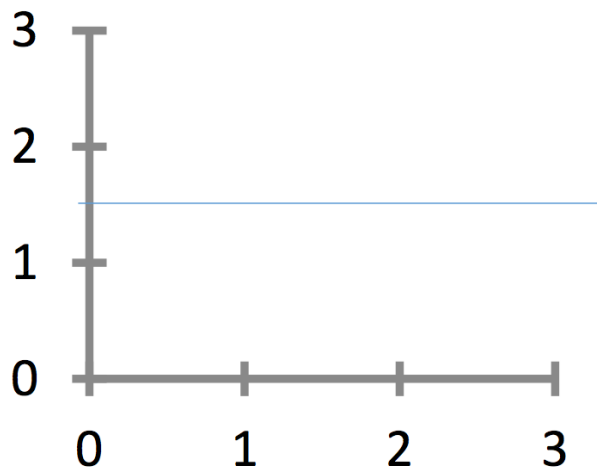Training Set

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

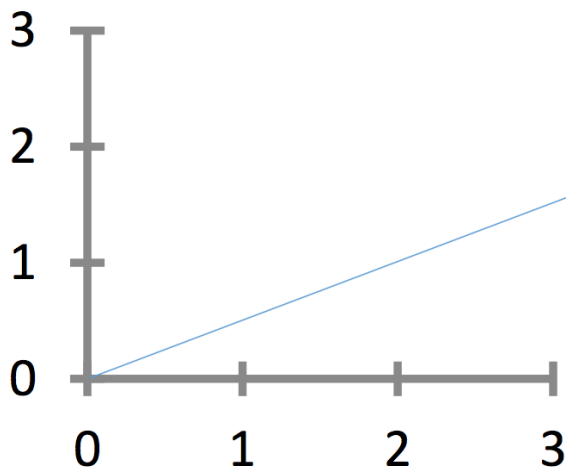Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

$\theta_i$'s:   Parameters

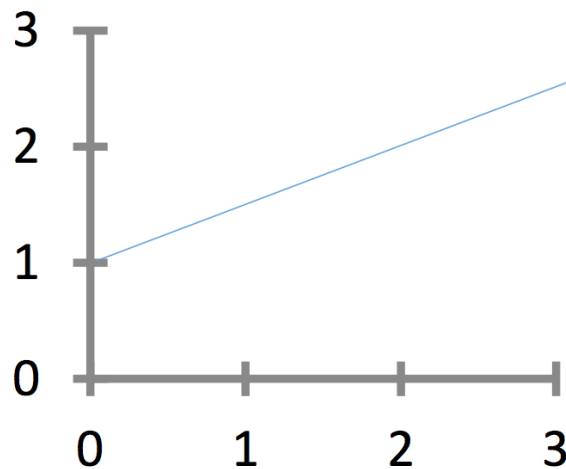How to choose $\theta_i$'s ?

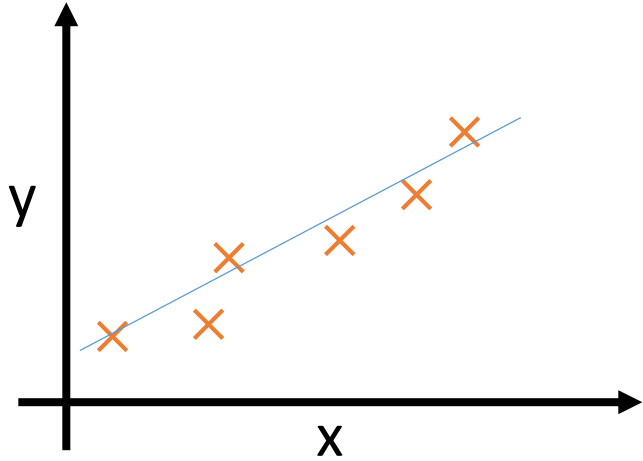$$h_\theta(x) = \theta_0 + \theta_1 x$$



$\theta_0 = 1.5$
$\theta_1 = 0$

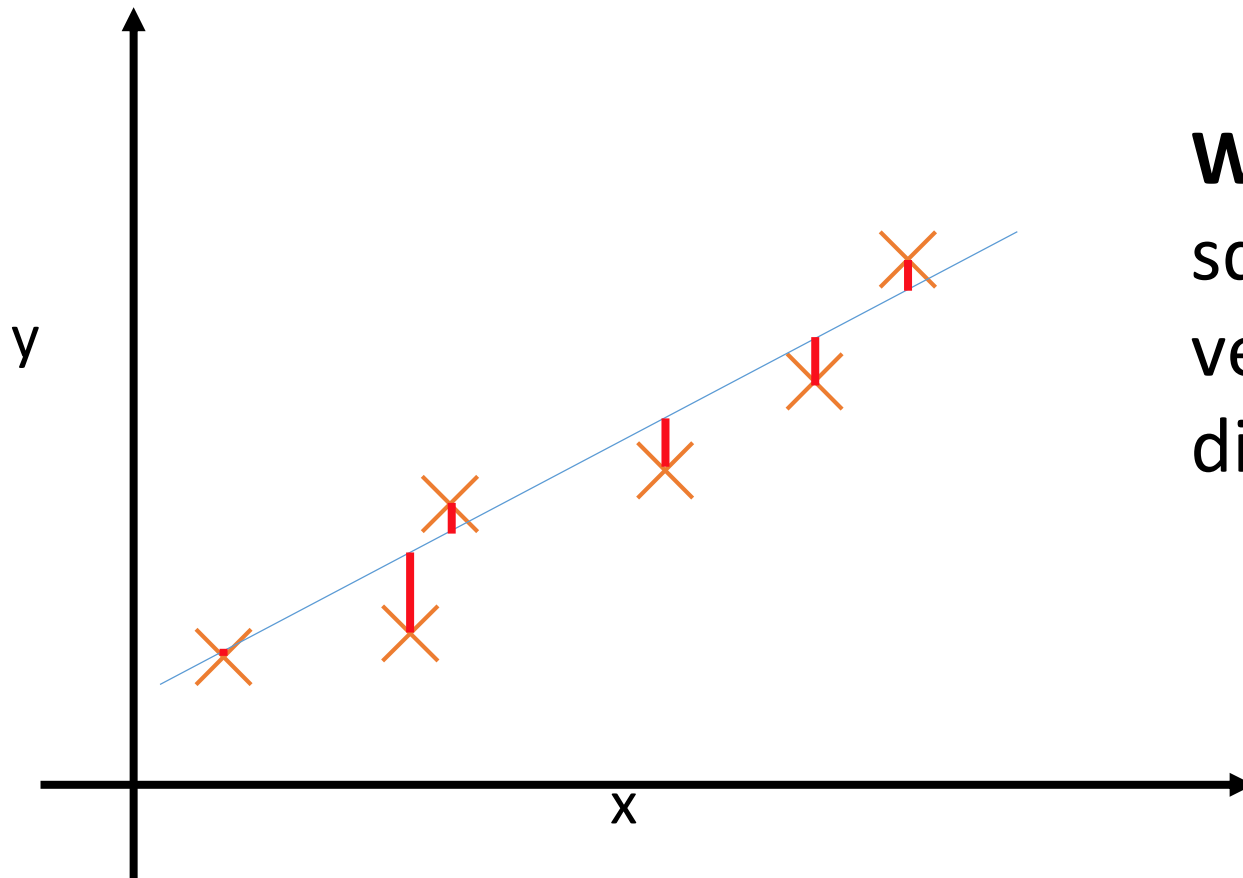$\theta_0 = 0$
$\theta_1 = 0.5$

$\theta_0 = 1$
$\theta_1 = 0.5$

**But what does "close" mean?**

Idea: Choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to $y$ for our training examples $(x, y)$

Quadratic cost function – on the board

**We choose:** squared vertical distance

Hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

Parameters: $\theta_0, \theta_1$

Cost Function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Goal: $\underset{\theta_0, \theta_1}{\text{minimize}} \; J(\theta_0, \theta_1)$

# Cost Function Surface Plot

# Contour Plots

- For a function F(x, y) of two variables, assigned different colours to different values of F
- Pick some values to plot
- The result will be *contours* – curves in the graph along which the values of F(x, y) are constant

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)



# $J(\theta_0, \theta_1)$

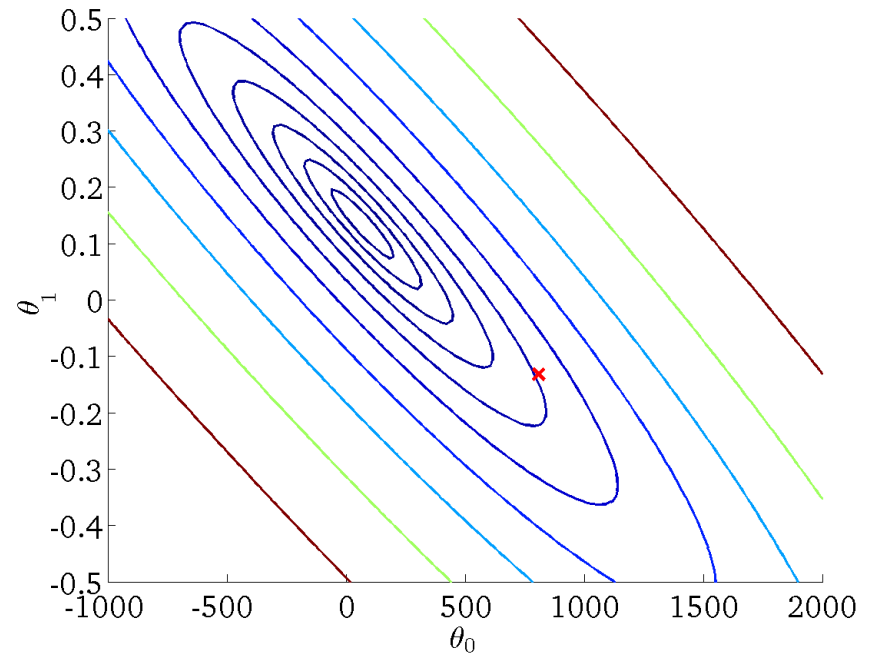(function of the parameters $\theta_0, \theta_1$)

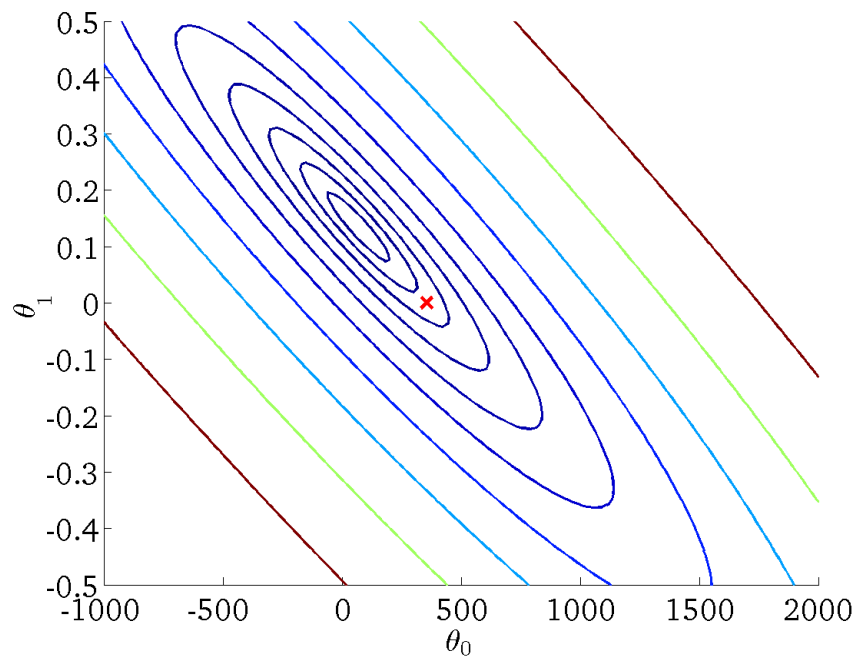# Cost Function Contour Plot

$$h_\theta(x)$$

(for fixed $\theta_0, \theta_1$ this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)

# $J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

Have some function $J(\theta_0, \theta_1)$

Want $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

**Outline:**

- Start with some $\theta_0, \theta_1$

- Keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$

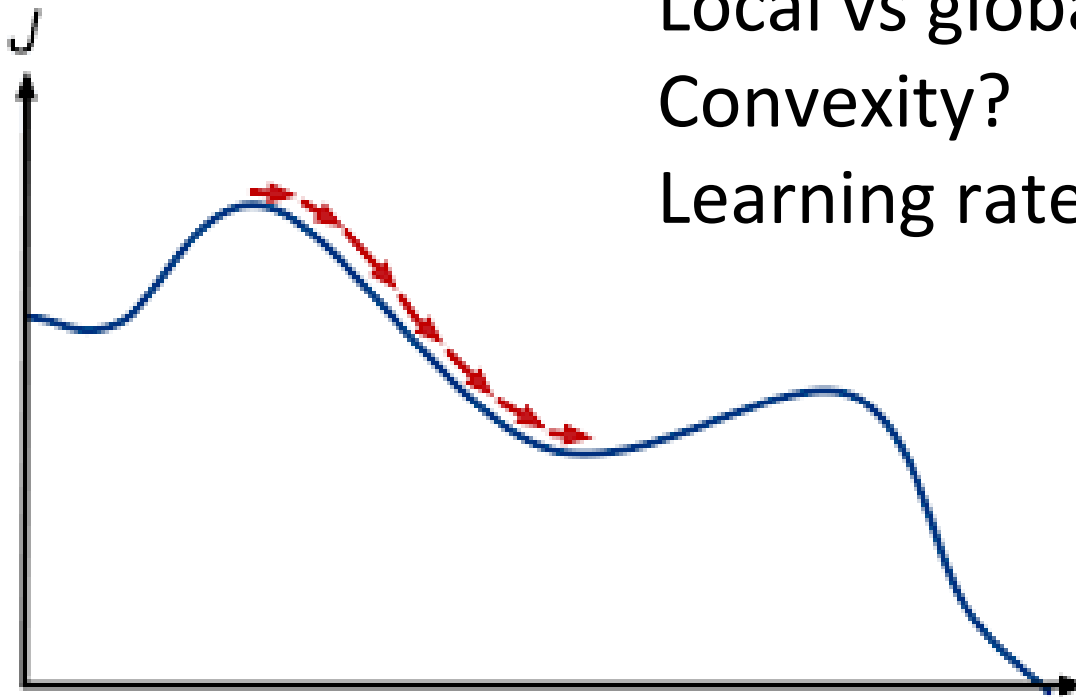    until we hopefully end up at a minimum

# Gradient Descent on the board

# Gradient Descent in 1D



To minimize f(x), we start with a random point and iterate with the update rule:

$$x_t \leftarrow x_{t-1} - \alpha \frac{df}{dx}(x_{t-1})$$

# Things to consider:



Local vs global minima?
Convexity?
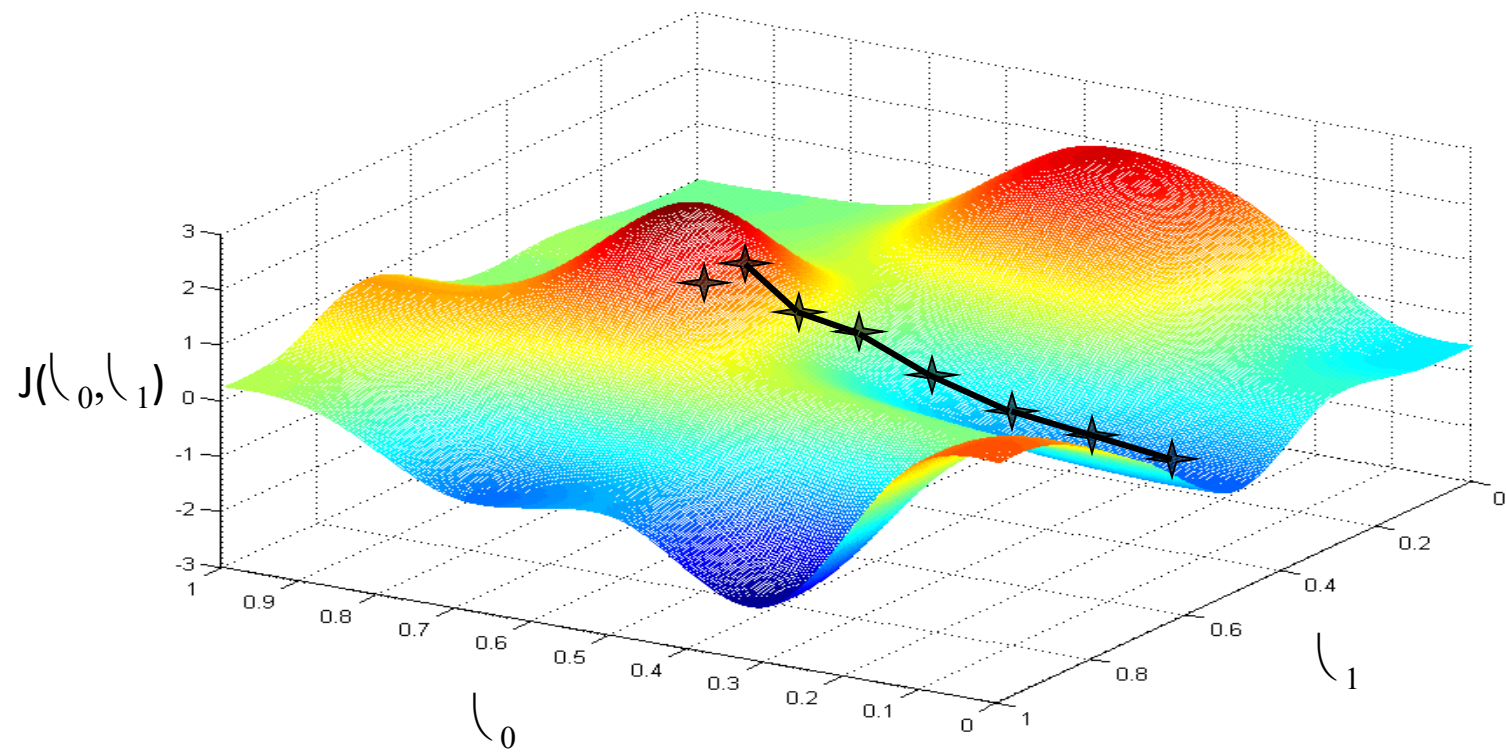Learning rate? ($\alpha$)

# Gradient Descent in Higher Dimensions



Update Rule:

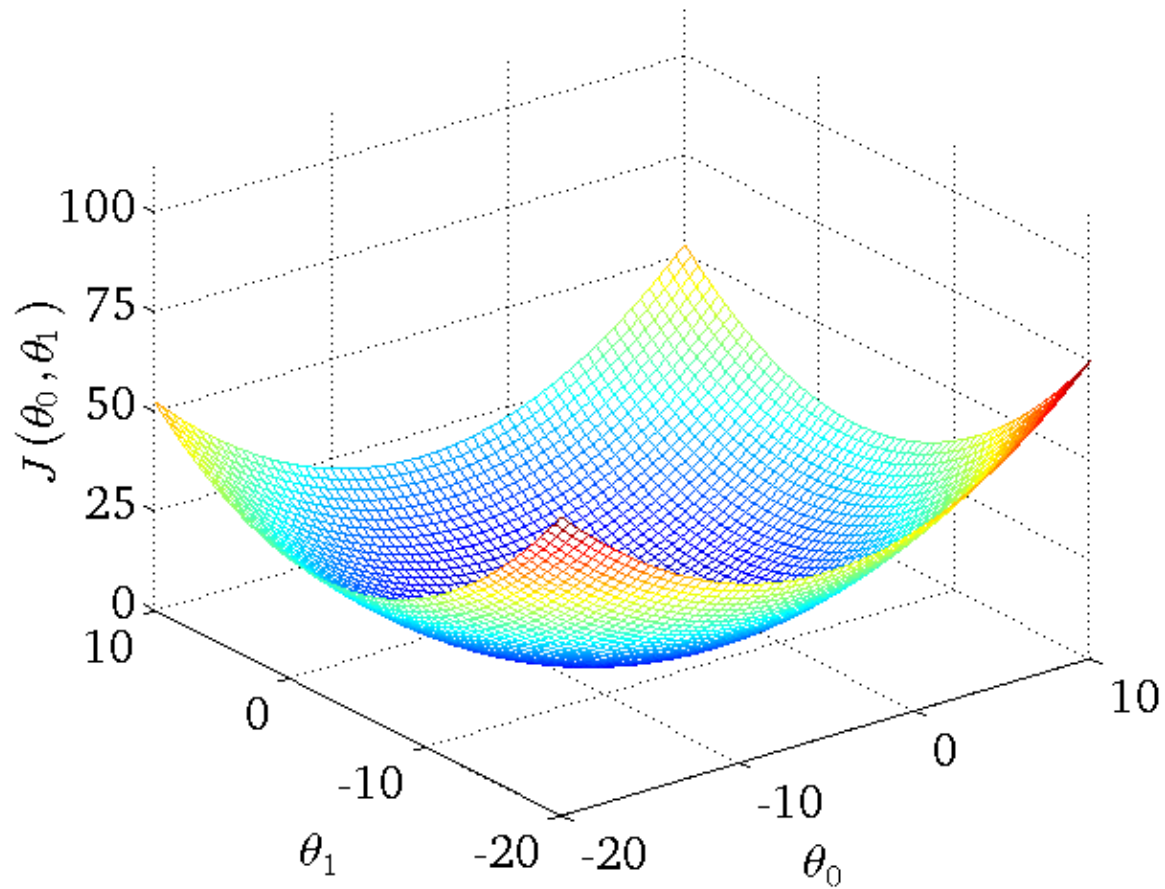$$x_t \leftarrow x_{t-1} - \alpha \nabla f(x_{t-1})$$

# Gradient, on the board

For Linear Regression, J is bowl-shaped ("convex")

# Gradient Descent Example
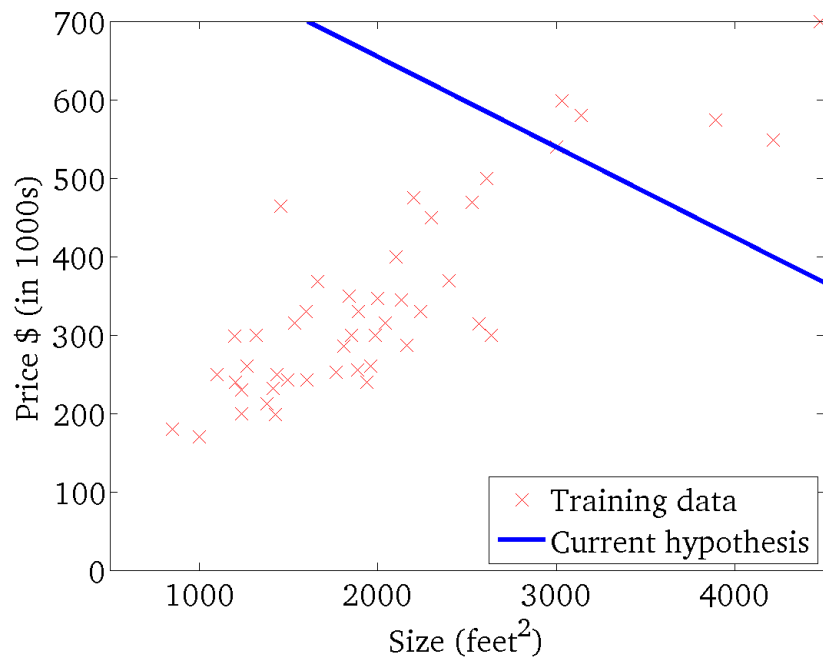
**Hypothesis:** $\quad h_\theta(x) = \theta_0 + \theta_1 x$

**Parameters:** $\quad \theta_0, \theta_1$

**Cost Function:** $\quad J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

**Goal:** $\quad \displaystyle\operatorname*{minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)



# $J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)

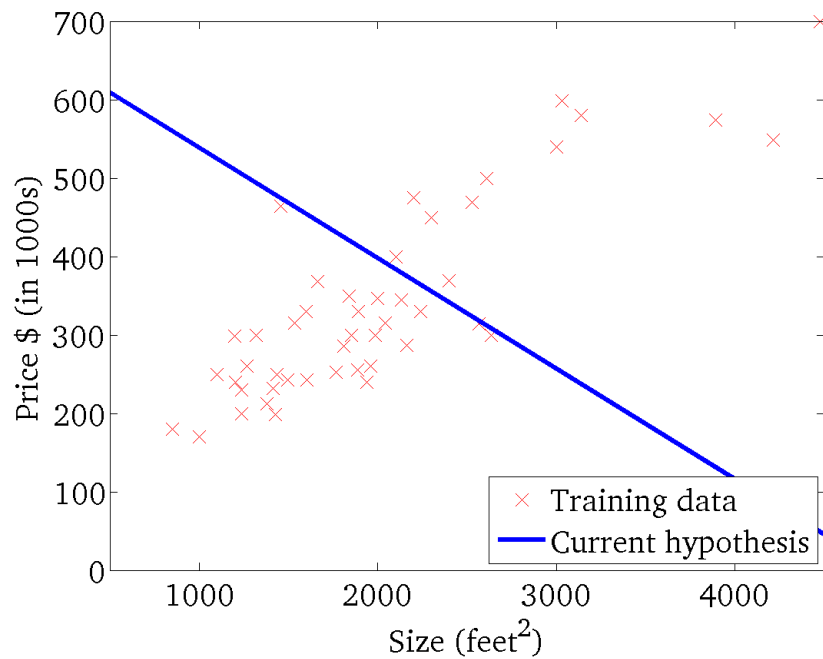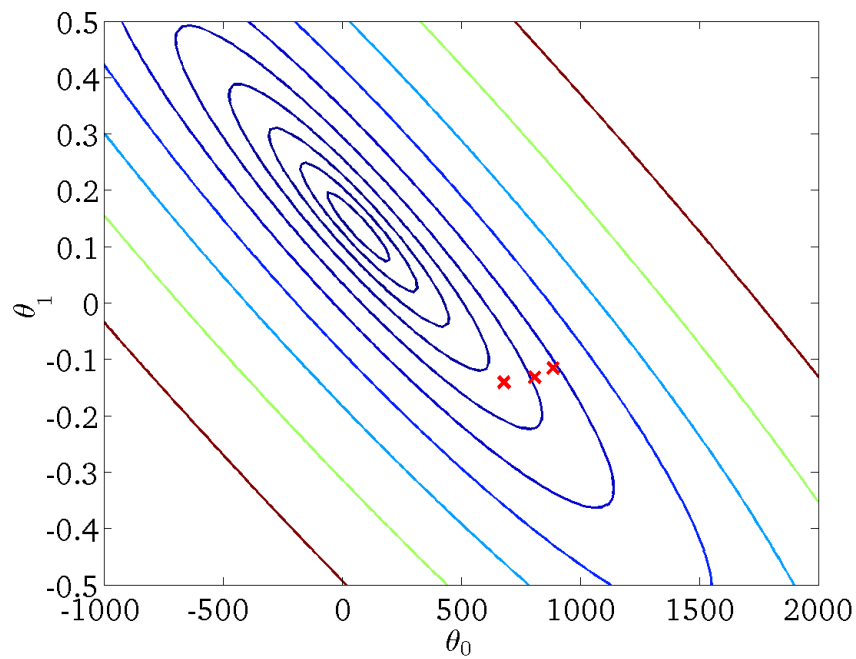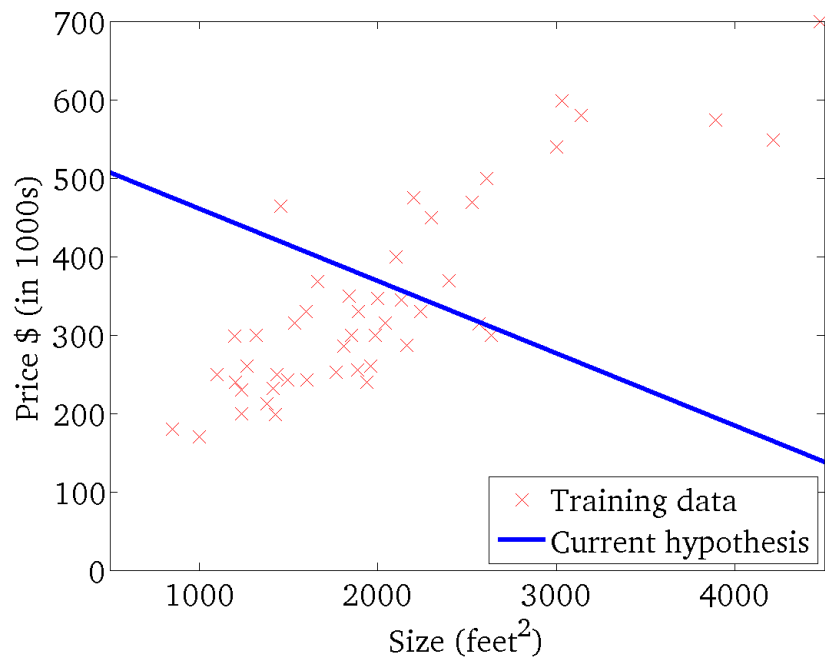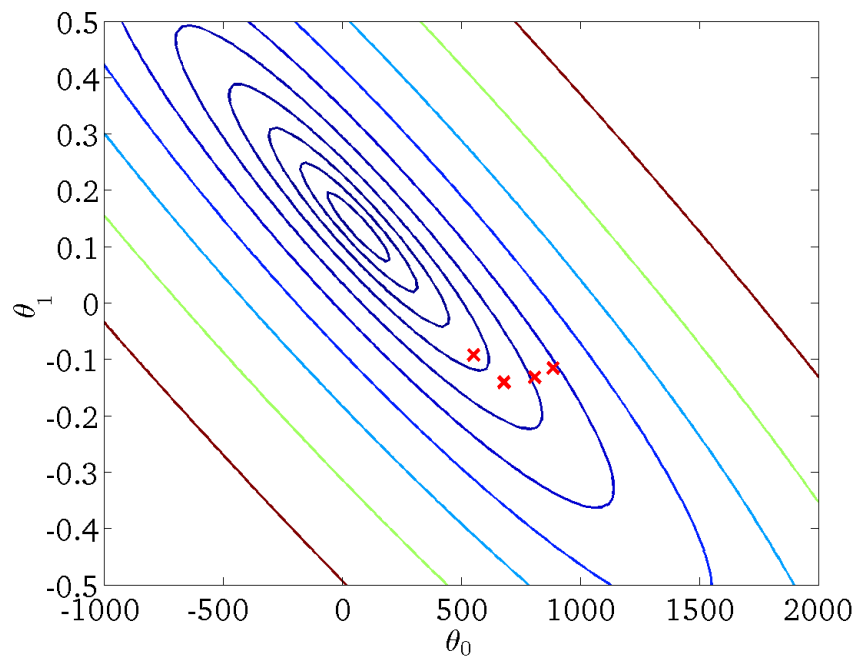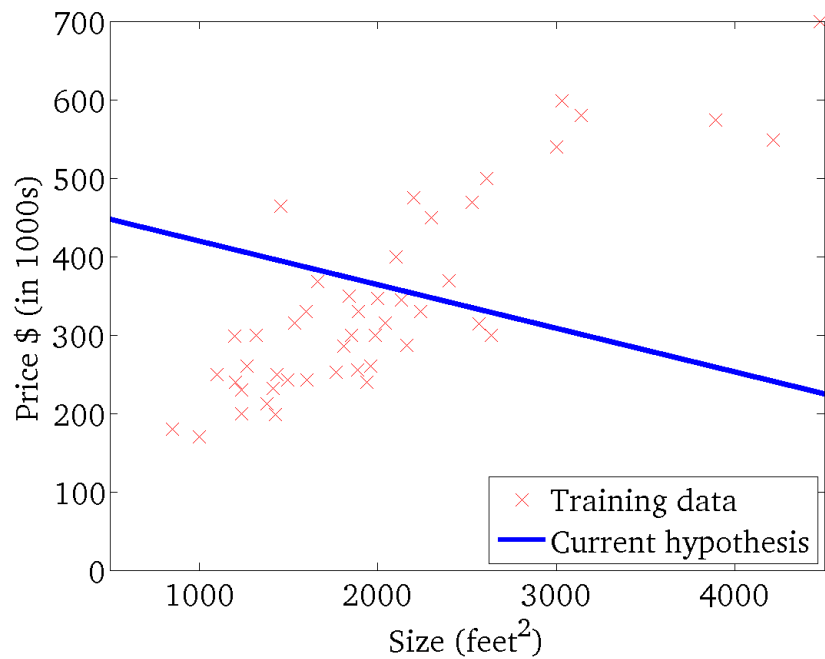# $J(\theta_0, \theta_1)$
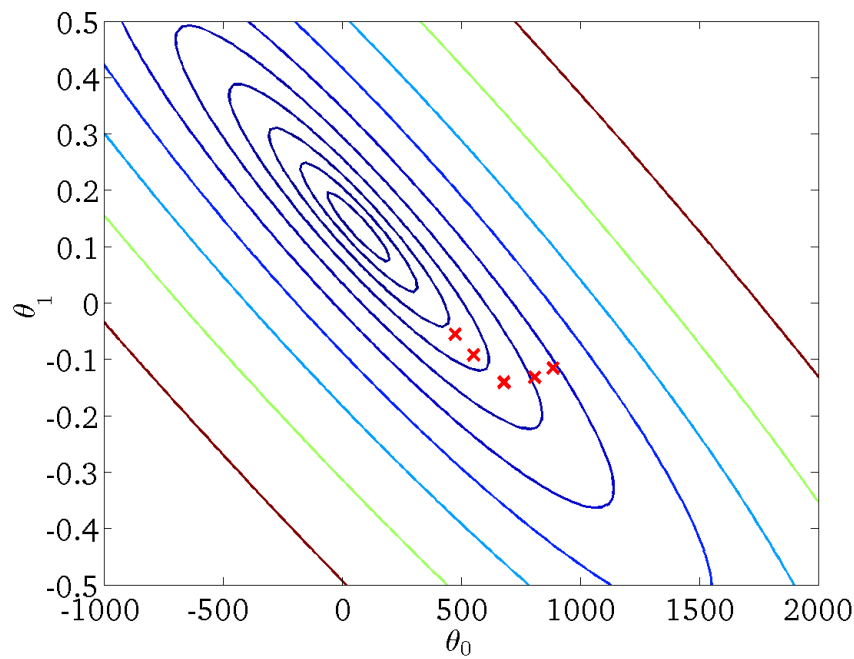
(function of the parameters $\theta_0, \theta_1$)

$h_\theta(x)$
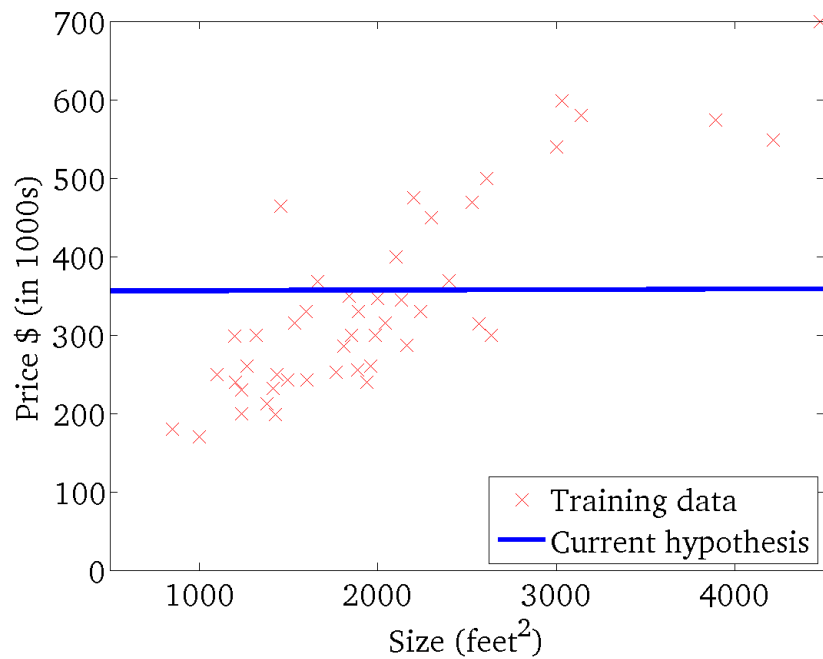
(for fixed $\theta_0, \theta_1$ this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

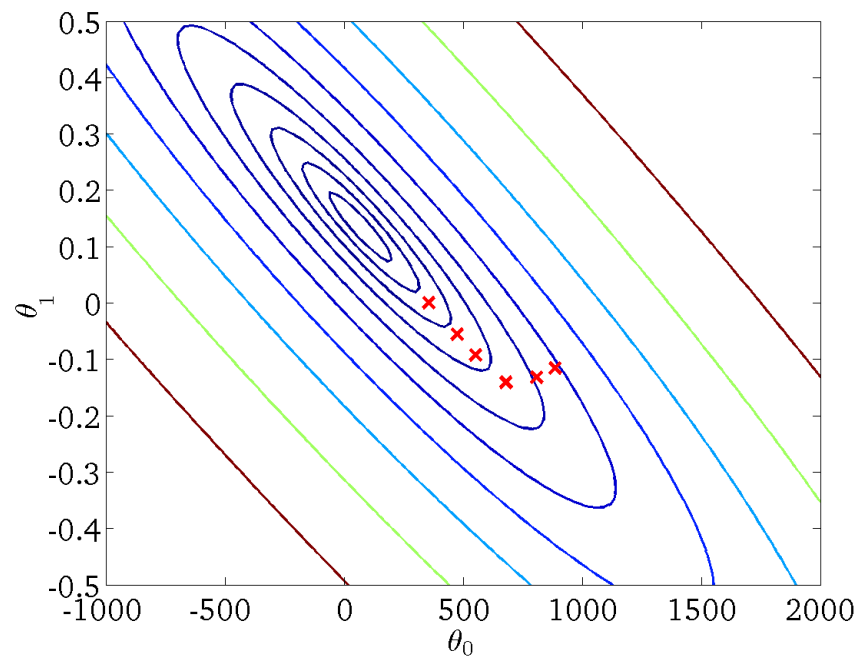# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)

# $J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)



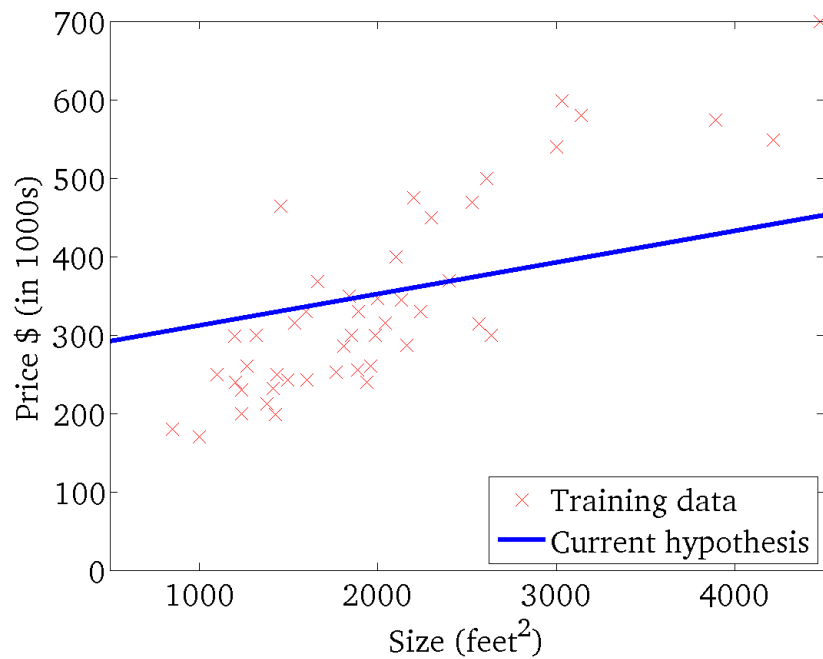# $J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)



# $J(\theta_0, \theta_1)$
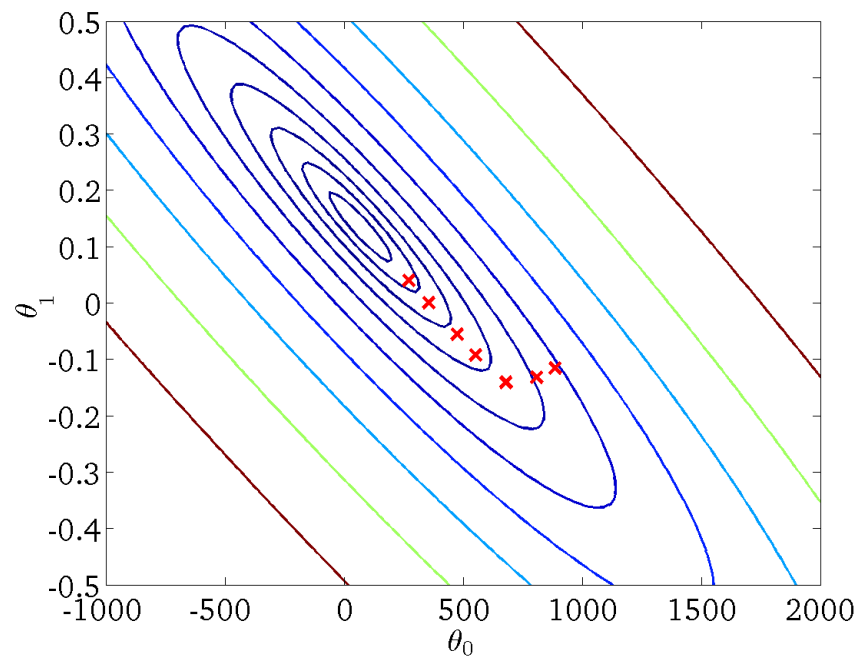
(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

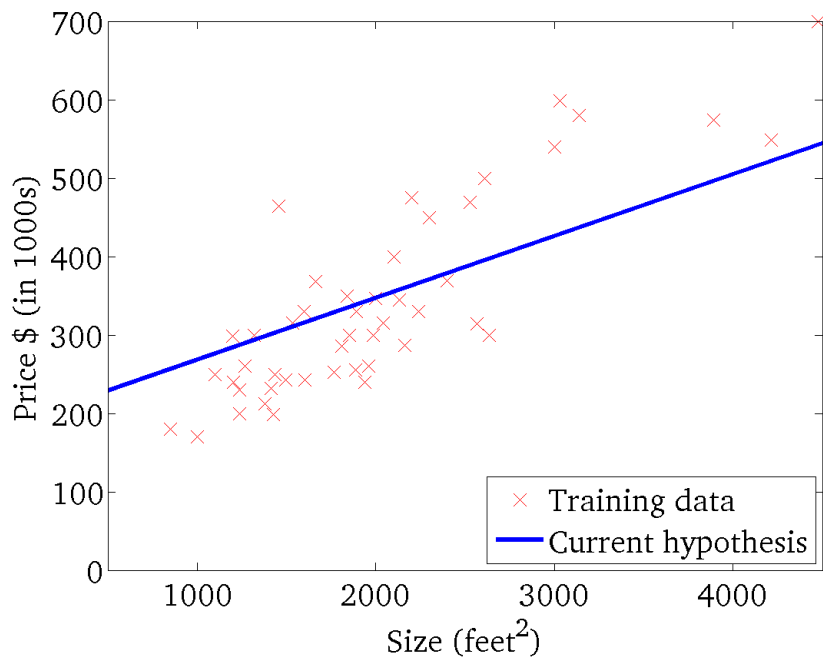(for fixed $\theta_0, \theta_1$ this is a function of x)

# $J(\theta_0, \theta_1)$
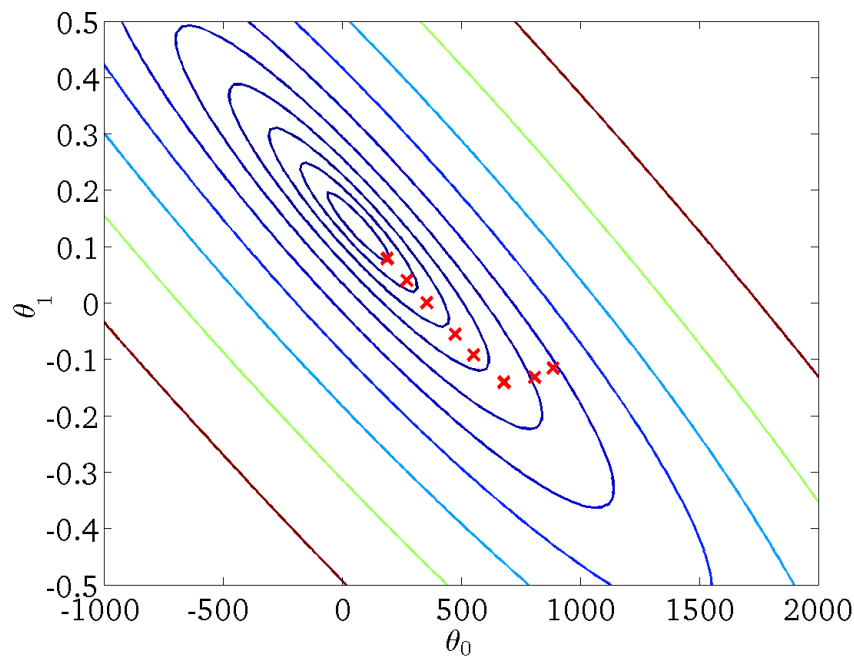
(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$ this is a function of x)



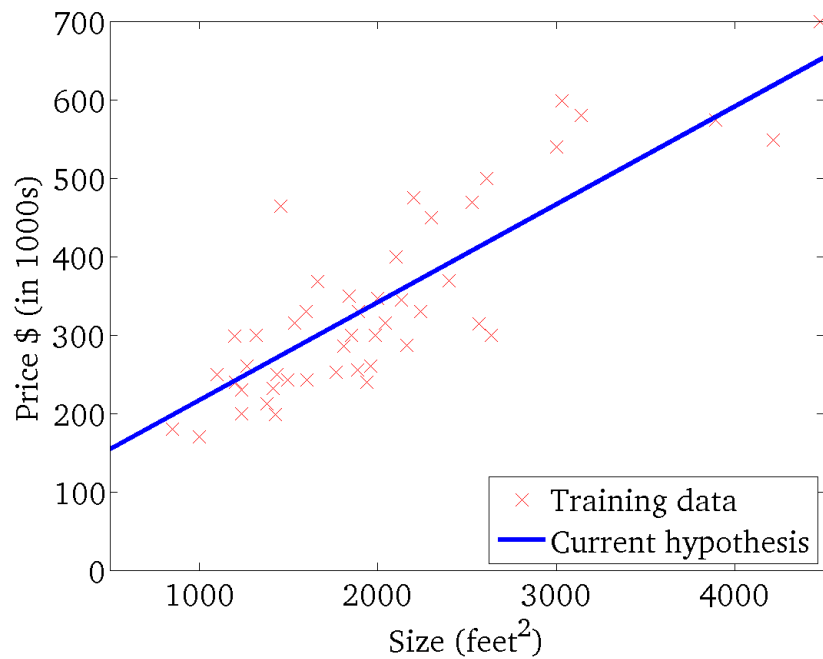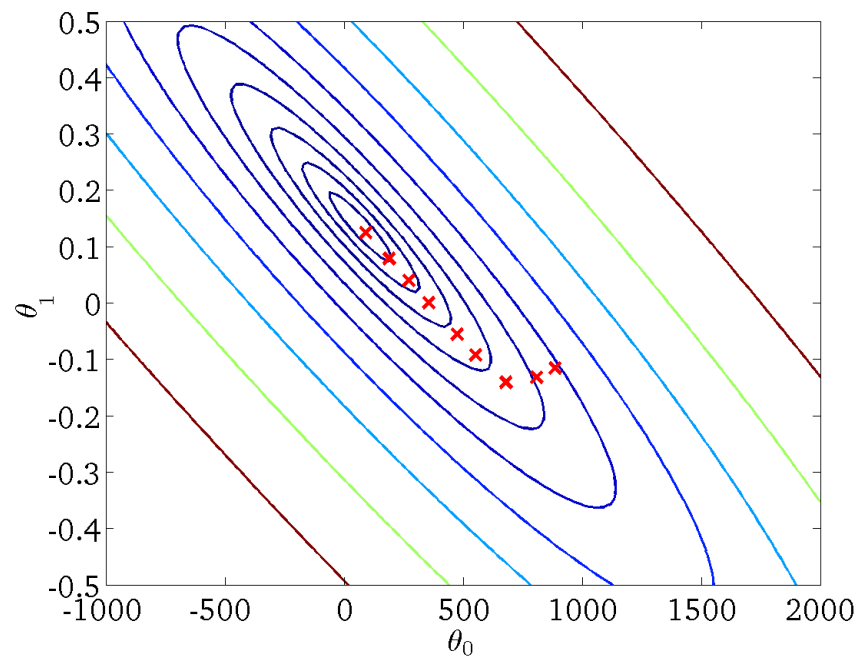# $J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

### (for fixed $\theta_0, \theta_1$ this is a function of x)
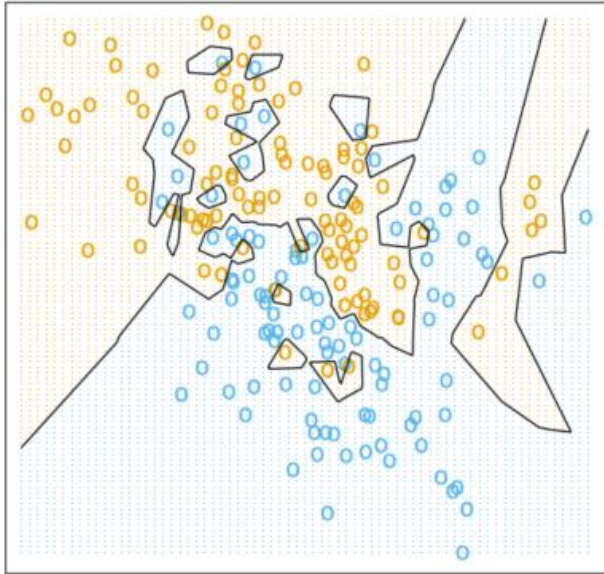


# $J(\theta_0, \theta_1)$

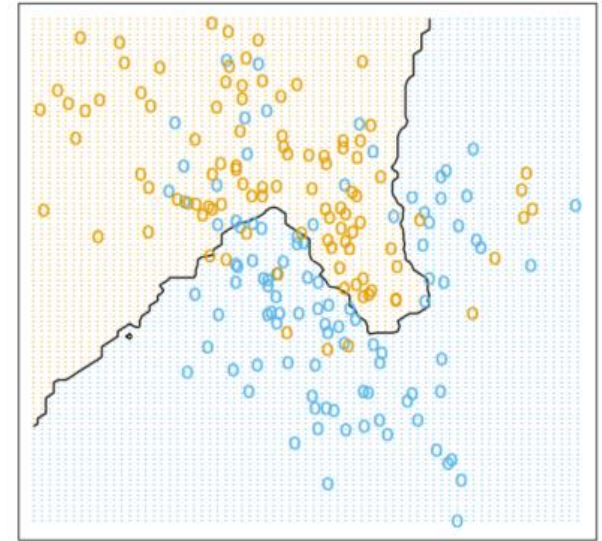### (function of the parameters $\theta_0, \theta_1$)

# Linear Regression vs. k-Nearest Neighbours
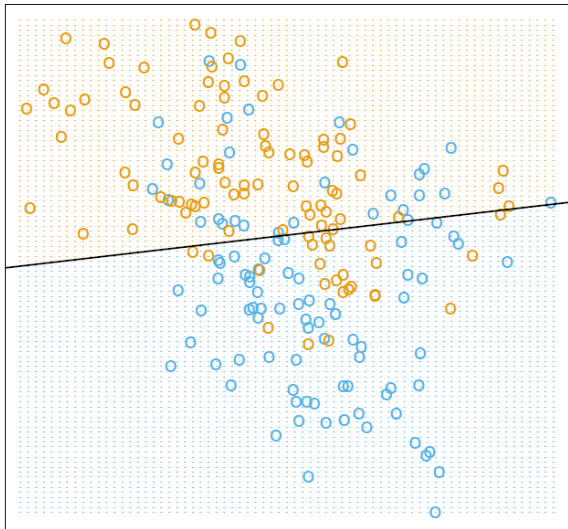


1-Nearest Neighbor Classifier

15-Nearest Neighbor Classifier

Linear Regression of 0/1 Response

Orange: y = 1
Blue: y = 0

# Linear Regression vs. k-Nearest Neighbours

- Linear Regression: the boundary can only be linear

- Nearest Neighbours: the boundary can more complex

- Which is better?

  - Depends on what the *actual boundary* looks like
  - Depends on whether we have enough data to figure out the *correct* complex boundary