

# CS M146: Introduction to Machine Learning

## Linear Regression

Aditya Grover

UCLA

The instructor gratefully acknowledges Sriram Sankararaman (UCLA) and Andrew Ng (Stanford) for some of the materials used in these slides, and many others who made their course materials freely available online.



<https://aditya-grover.github.io/>



@adityagrover\_

# Recap: Overview of ML

**Definition:** Machine Learning (ML) is the study of algorithms that improve their **performance P**, at some **task T**, with **experience E**.  
A well-defined learning task is given by  $\langle T, P, E \rangle$ .  
[Tom Mitchell (1998)]

## 3 Types of Learning:

Supervised, Unsupervised, Reinforcement Learning

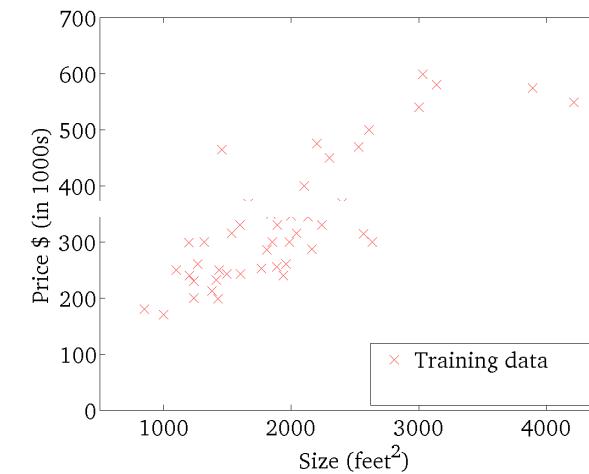
# Plan for Today

- Supervised Learning
  - **Linear Regression**

# Regression: Example

## Housing prices in Portland

Living area (feet <sup>2</sup> )	Price (1000\$)
2104	400
1600	330
2400	369
1416	232
3000	540
:	:



Given a dataset of housing prices, can we learn to predict the prices of other houses in Portland as a function of their living area sizes?

# Regression – Problem Setup

## Given

A training dataset consisting of  $n$  labelled **training instances**

- Input **features/attributes**  $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^d$
- Corresponding **labels**  $\mathbf{y} = \{y^{(1)}, \dots, y^{(n)}\}$  where  $y^{(i)} \in \mathbb{R}$

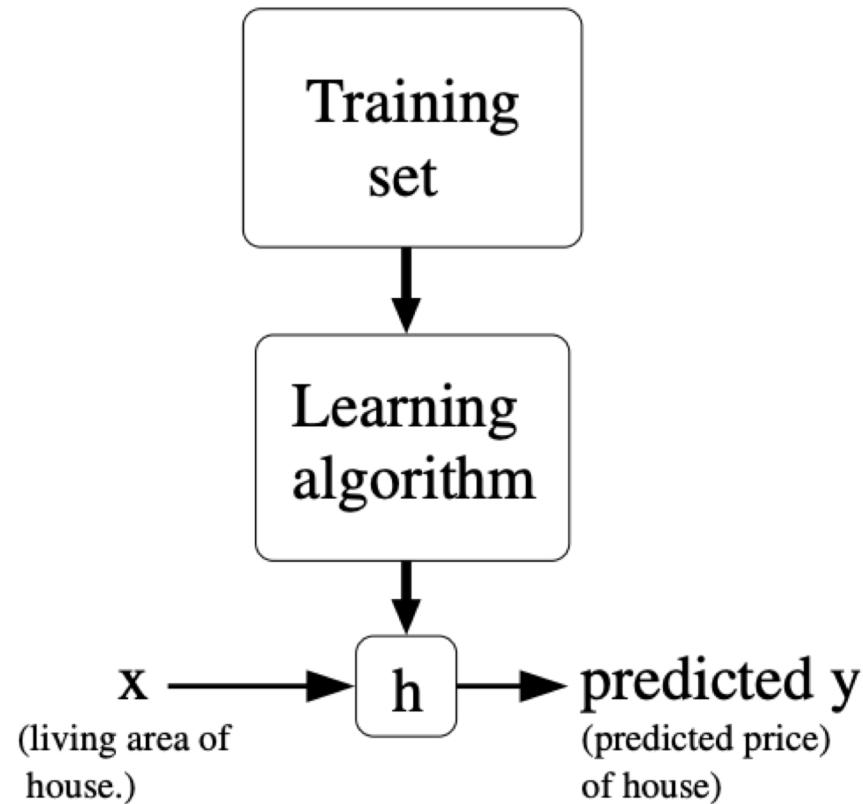
## Output

A **hypothesis function**  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  from a **hypothesis class**  $h \in \mathcal{H}$  such that  $h(\mathbf{x}) \approx y$

# Aside: Notation scheme for the course

- Every time we see a new term, it will be colored **blue**
- Vectors are always in **bold**, e.g.,  $\mathbf{x}^{(1)}, \theta$
- Vector elements and scalars are not in bold, e.g.,  $x_1^{(1)}, \theta_2, y^{(3)}$
- Why  $\mathbf{x}^{(1)}$  as opposed to:
  - $x^1$  to avoid confusion with exponentiation
  - $x_1$  to avoid confusion with indexing of vector elements
- **If you spot a typo, let us know and we'll fix it!**

# Regression – 3 Step Recipe



Steps for learning  $h$

1. Represent hypothesis class  $\mathcal{H}$
2. Define a **loss**
3. Optimize loss to find best  $h \in \mathcal{H}$

# Linear Regression – Step 1

- In **linear** regression, we represent hypothesis class  $\mathcal{H}$  as the set of all **parameterized linear functions**

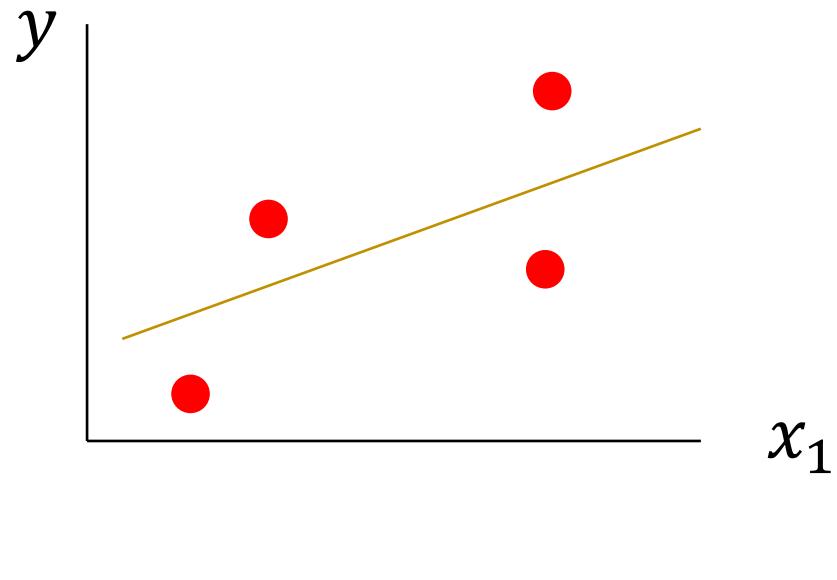
$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_d x_d$$

- **Note:** linearity is defined with respect to  $\theta_i$ 's (and not  $x_i$ 's)
- $\boldsymbol{\theta} = [\theta_0, \theta_1, \dots, \theta_d] \in \mathbb{R}^{d+1}$  are the **parameters** of hypothesis class  $\mathcal{H}$ .
- Alternate terminologies:
  - $\theta_0$  : **bias**
  - $\theta_{1:d}$  : **weights**
- **Vectorization:** Define  $x_0 = 1$  such that  $\mathbf{x} \in \mathbb{R}^{d+1}$

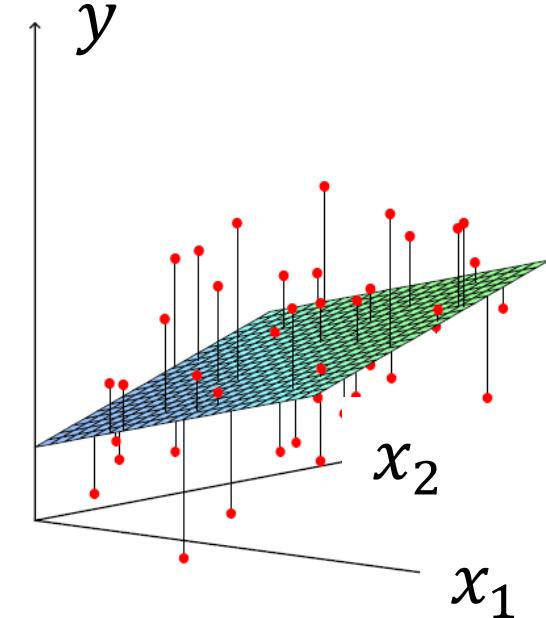
$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j = \boldsymbol{\theta}^T \mathbf{x}$$

# Geometric Intuition

$d = 1, h_{\theta}(x)$ : straight line



$d = 2, h_{\theta}(x)$ : plane

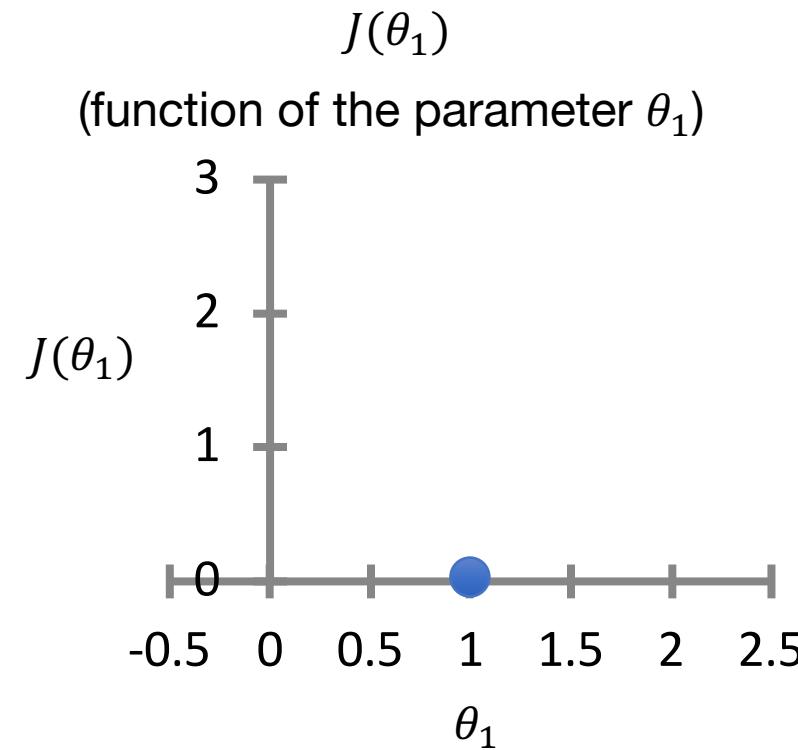
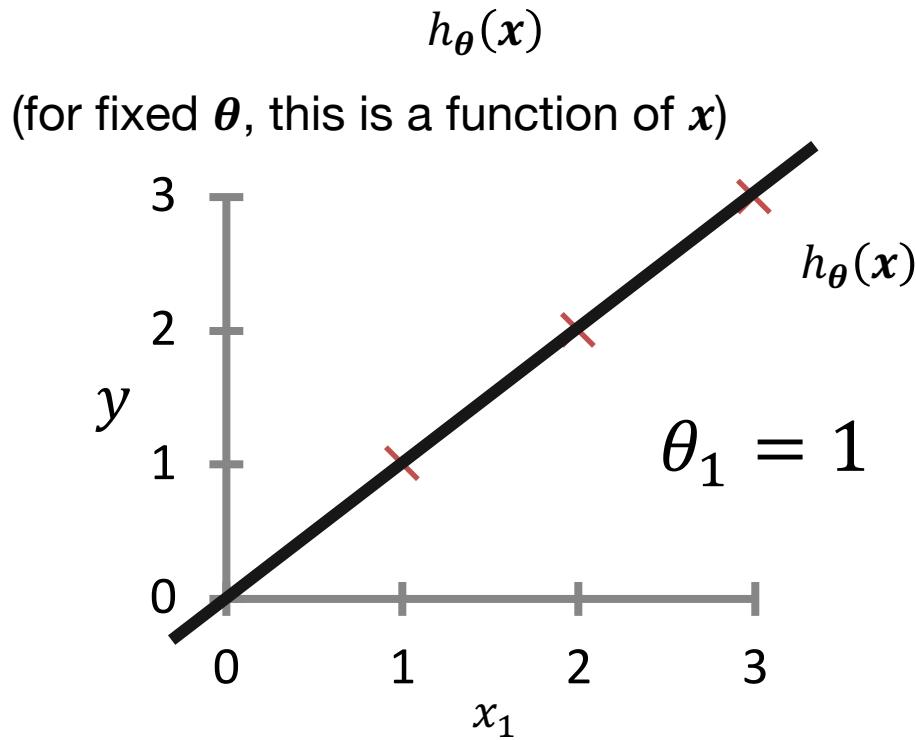


$d > 2, h_{\theta}(x)$ : hyperplane

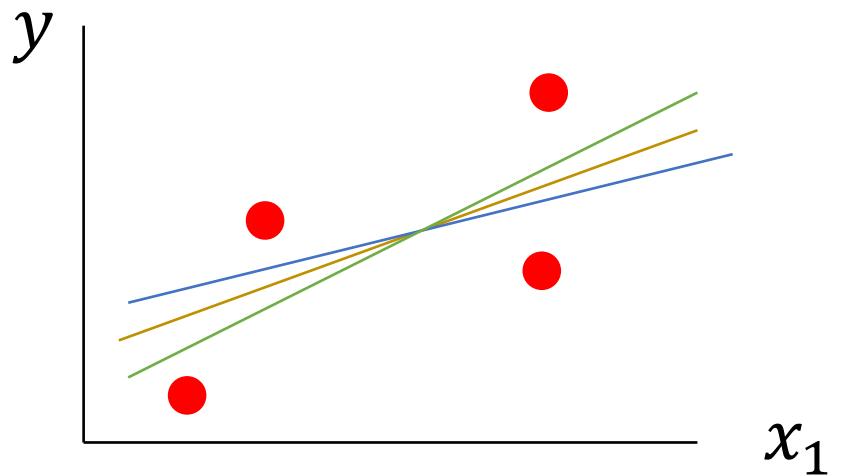
# Intuition Behind Loss Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

For intuition, let us assume  $d = 1$  and  $\theta_0 = 0$ . Hence  $\boldsymbol{\theta} = [0, \theta_1]$ .



# Linear Regression – What Next?



Steps for learning  $h$

1. Represent hypothesis class  $\mathcal{H}$
2. **Define a loss objective**
3. Minimize loss to find best  $h \in \mathcal{H}$

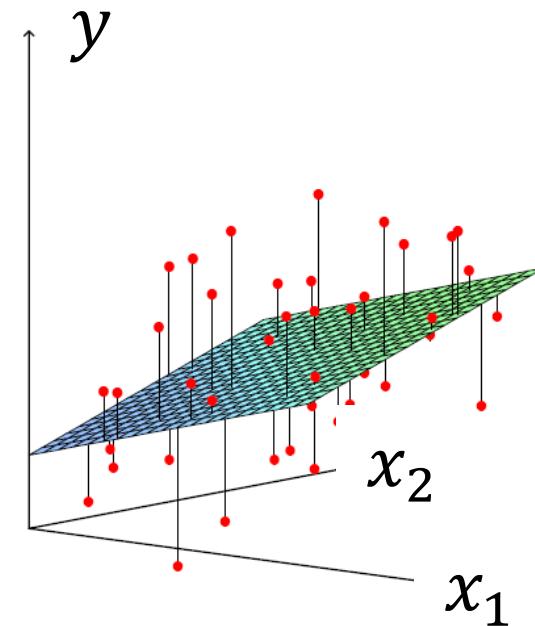
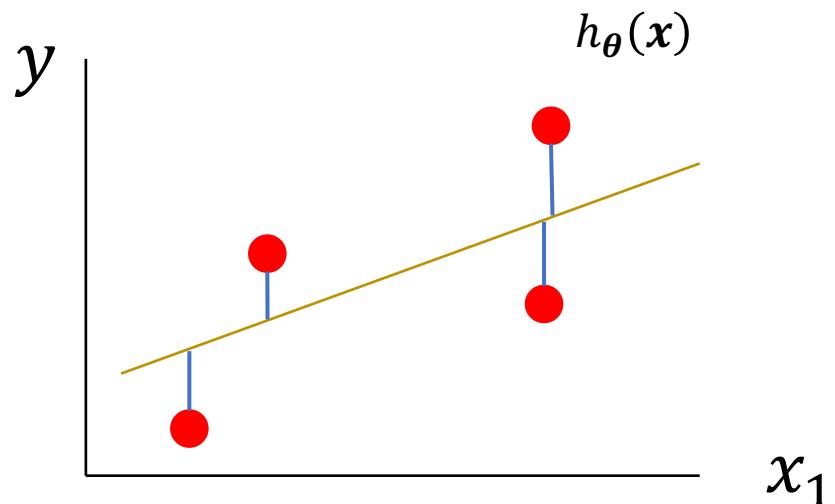
Many hypothesis in  $\mathcal{H}$ .

Which one is best hypothesis? – Step 2

# Least Squares Linear Regression – Step 2

- **Loss function:** A way to score hypothesis. Lower loss is better.
- Least squares loss function

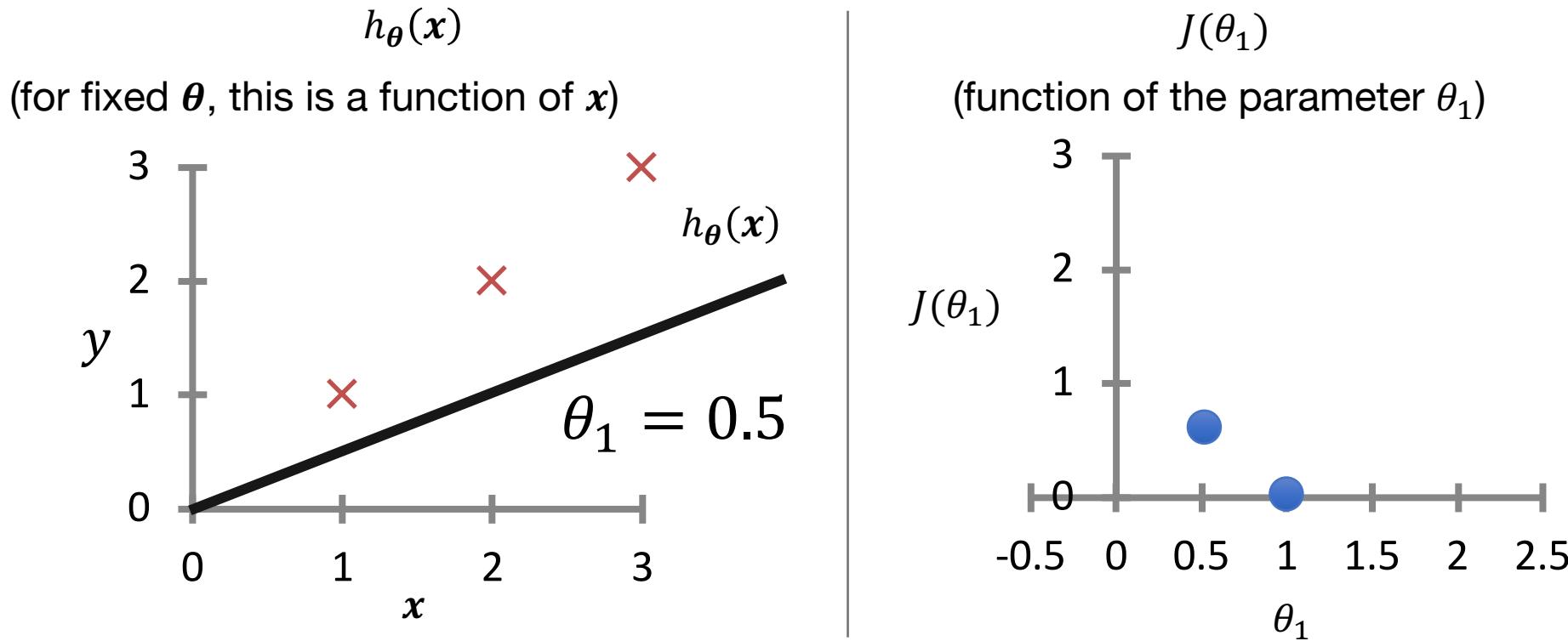
$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$$



# Intuition Behind Loss Function

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

For intuition, let us assume  $d = 1$  and  $\theta_0 = 0$ . Hence  $x \in \mathbb{R}$  so  $\boldsymbol{\theta} = [0, \theta_1]$ .

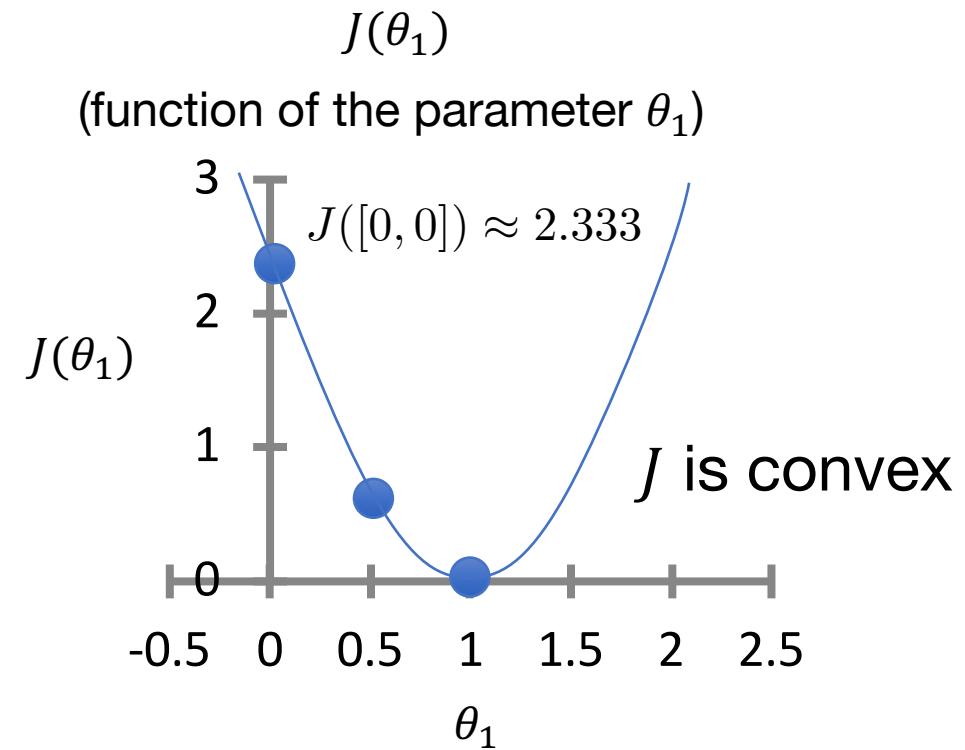
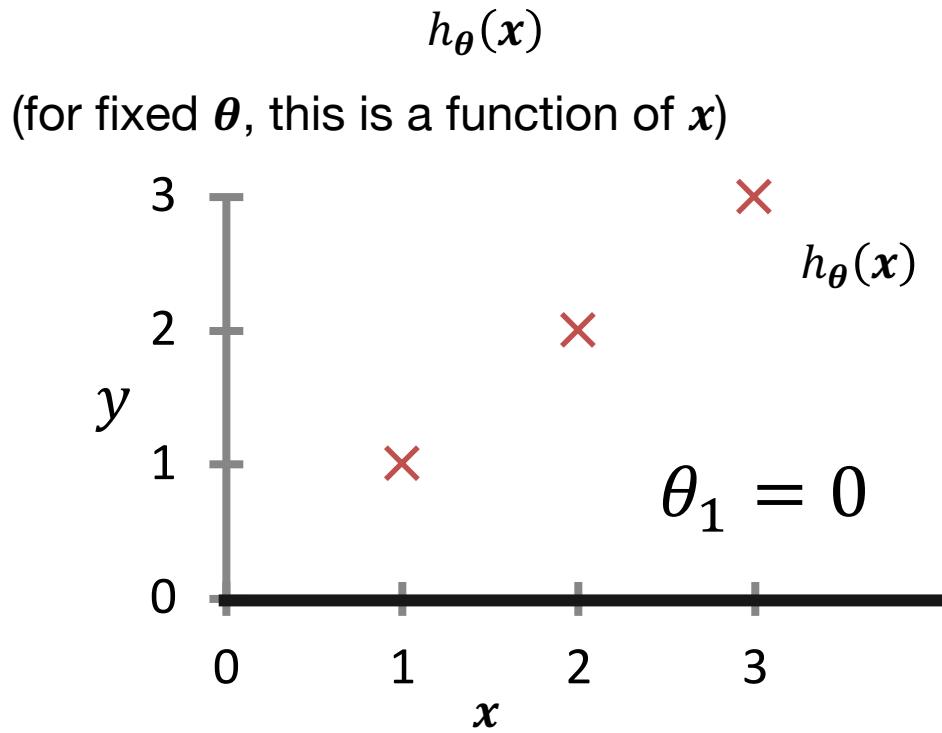


$$J([0, 0.5]) = \frac{1}{2 \times 3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \approx 0.58$$

# Intuition Behind Loss Function

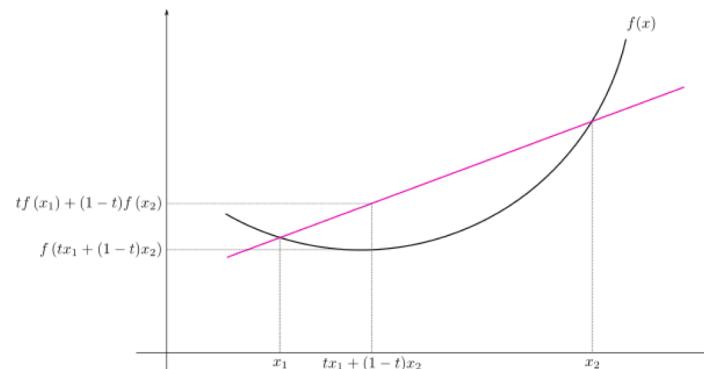
$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

For intuition, let us assume  $d = 1$  and  $\theta_0 = 0$ . Hence  $x \in \mathbb{R}$  so  $\boldsymbol{\theta} = [0, \theta_1]$ .



# Background: Convex Functions

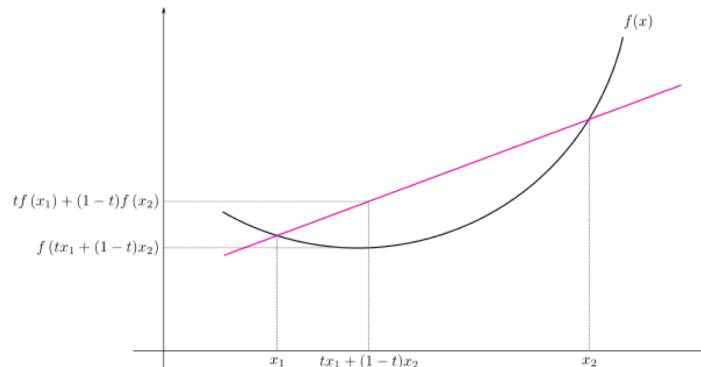
- **Definition:** A real-valued function  $f$  is **convex** if the line segment joining any two points lies above the function



- A function  $f$  is convex iff for all  $0 \leq t \leq 1$  and all  $x_1, x_2 \in X$ :  
$$f(tx_1 + (1 - t)x_2) \leq tf(x_1) + (1 - t)f(x_2)$$
- For **strictly convex** functions, replace  $\leq$  by  $<$  above

# Background: Convex Functions

- **(Stated without proof)** Convex functions have a unique optima

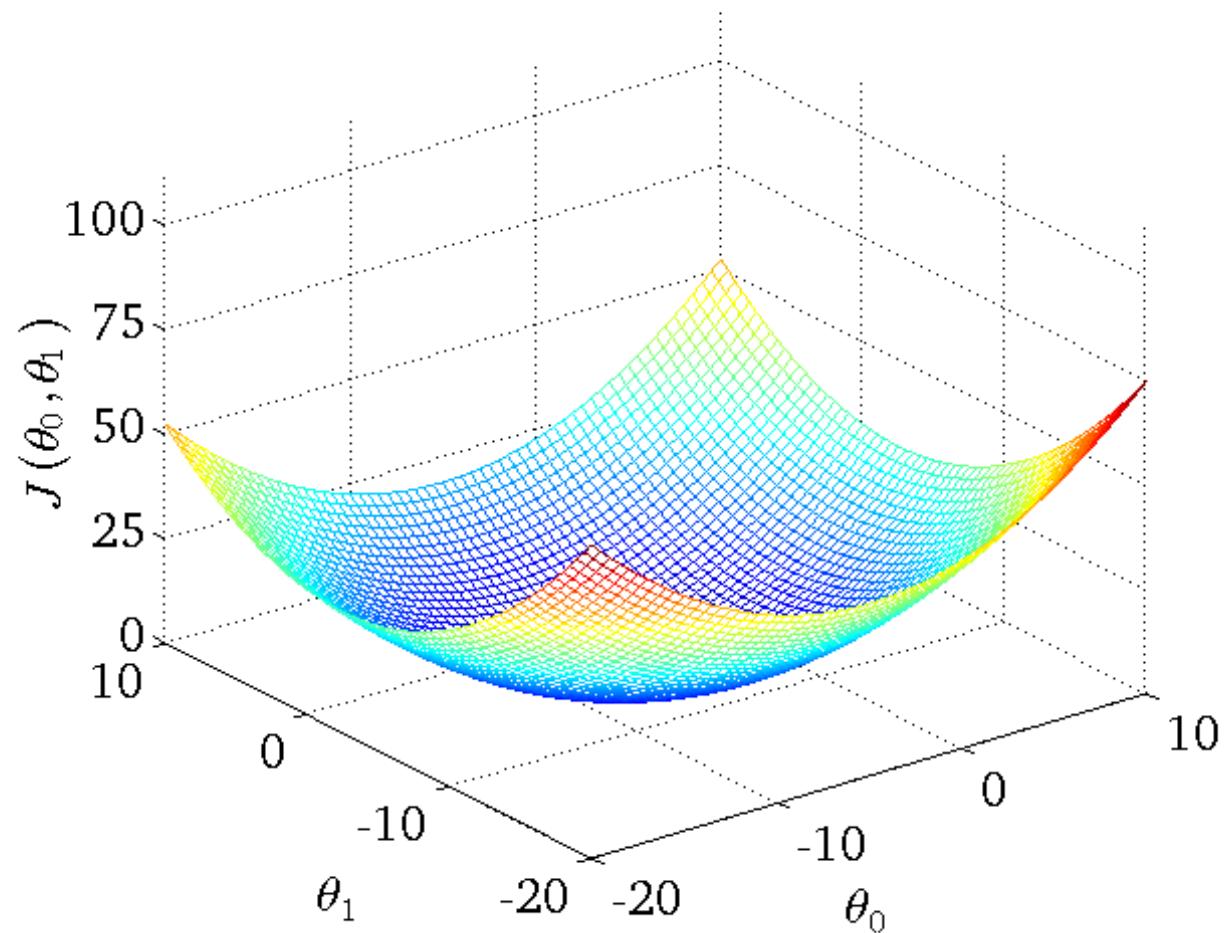


- **(Stated without proof)**

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$$

$J(\theta)$  is **convex** and **differentiable** w.r.t.  $\theta$

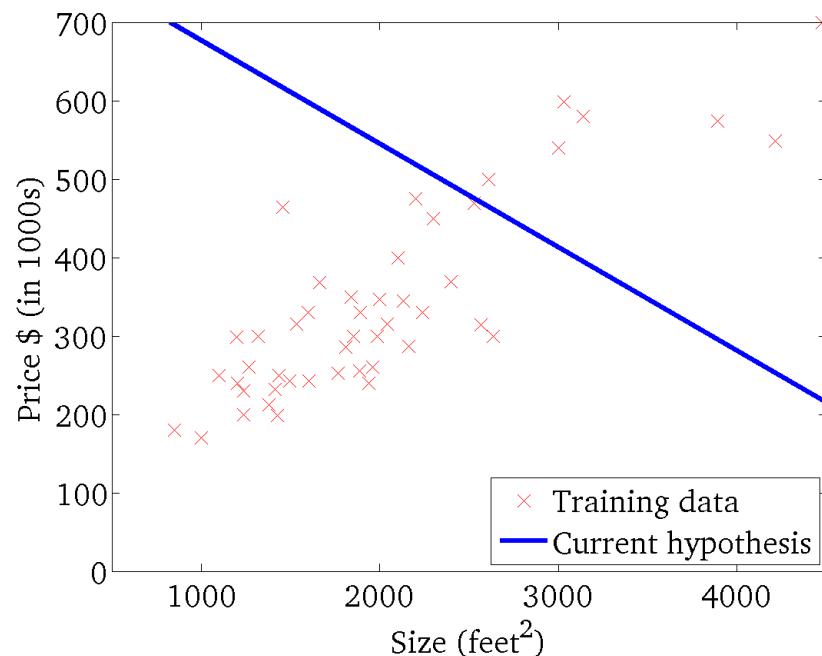
# Intuition Behind Loss Function



# Intuition Behind Loss Function

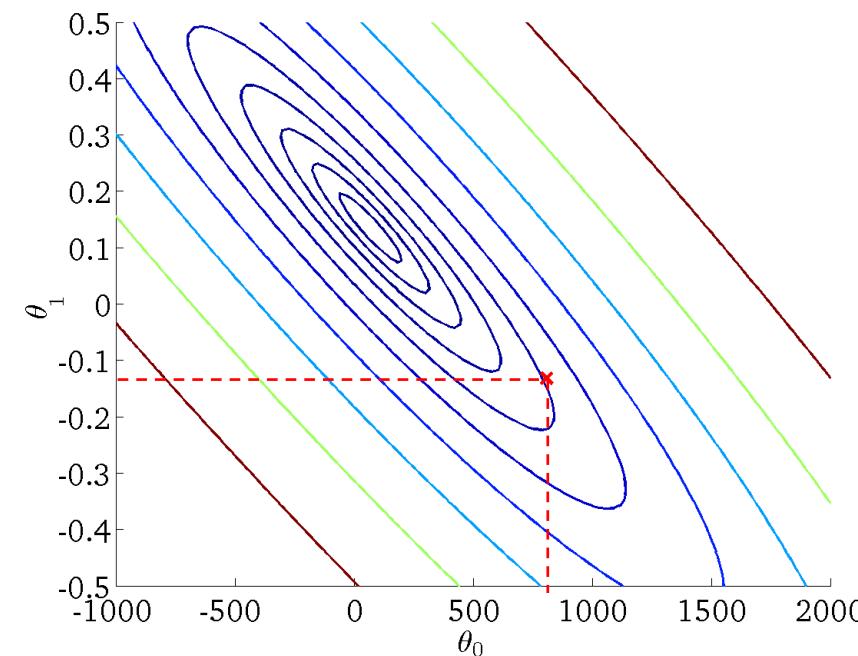
$$h_{\theta}(x)$$

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



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

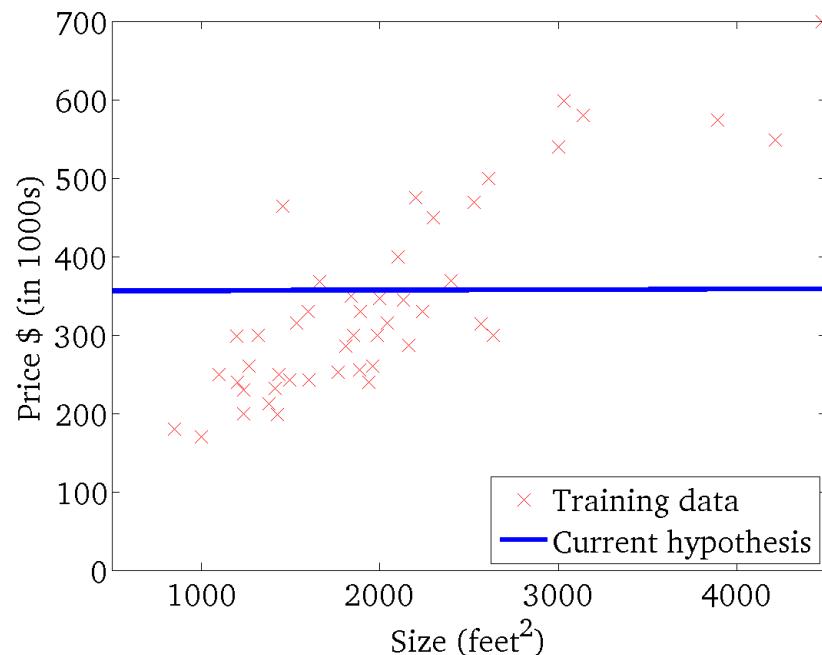
(function of the parameters  $\theta_0, \theta_1$ )



# Intuition Behind Loss Function

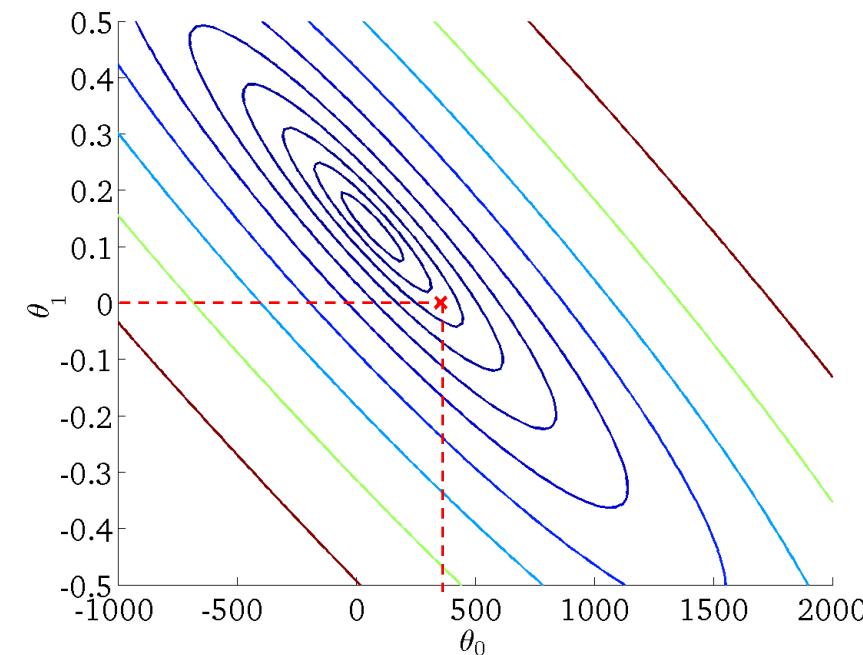
$$h_{\theta}(x)$$

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



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

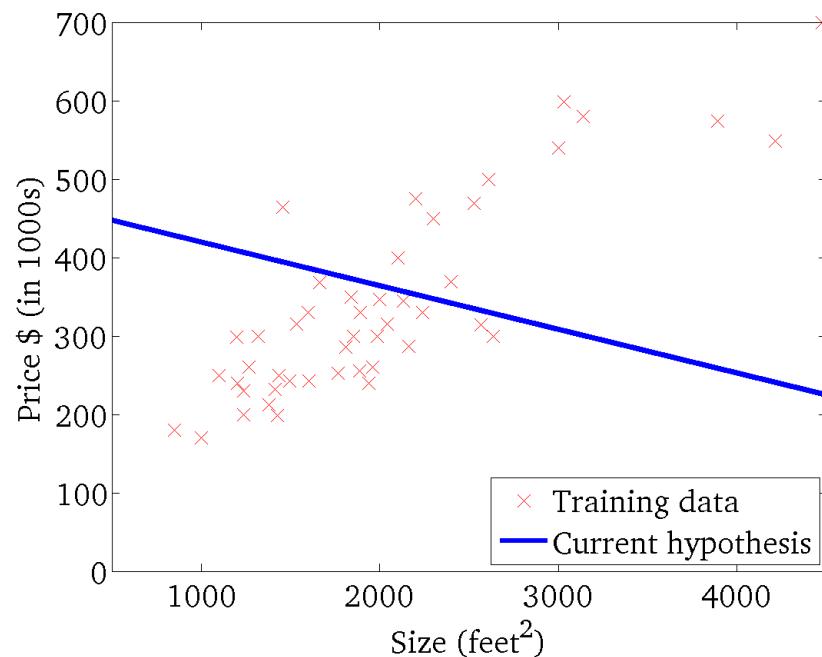
(function of the parameters  $\theta_0, \theta_1$ )



# Intuition Behind Loss Function

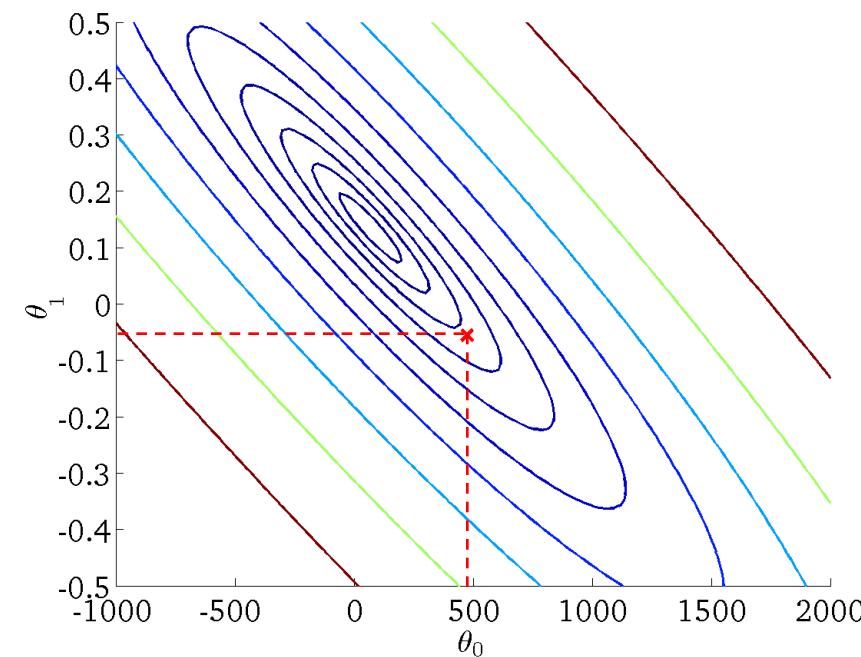
$$h_{\theta}(x)$$

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



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

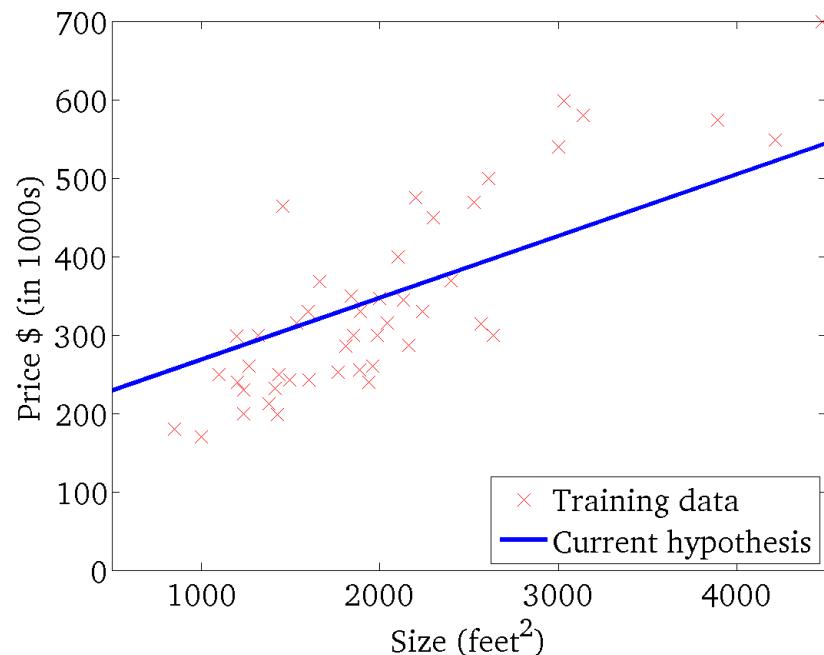
(function of the parameters  $\theta_0, \theta_1$ )



# Intuition Behind Loss Function

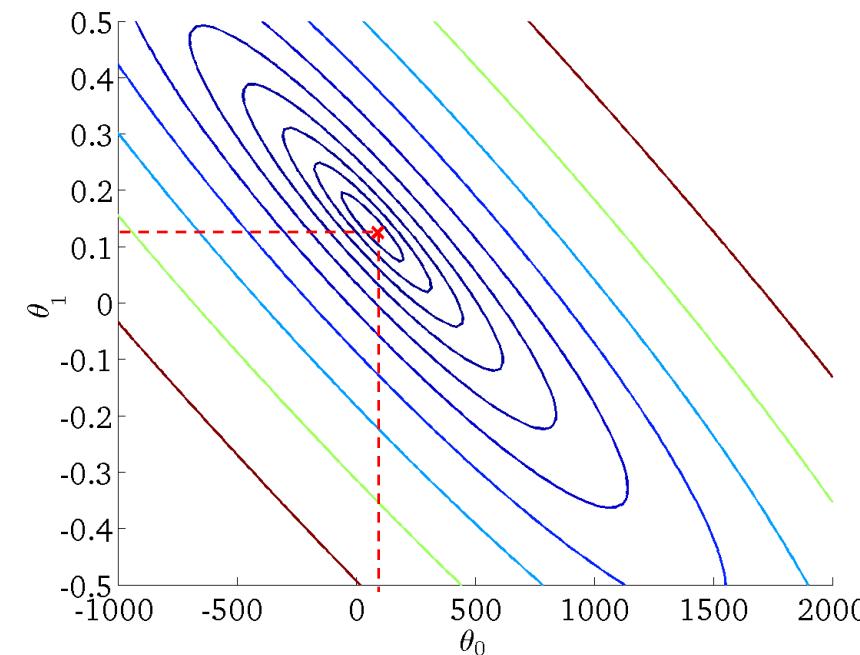
$$h_{\theta}(x)$$

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

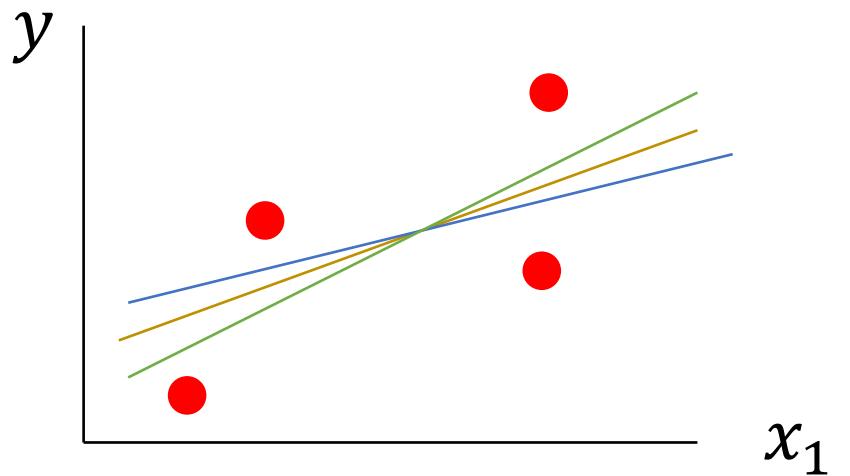


$$J(\theta_0, \theta_1)$$

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



# Linear Regression – What Next?



Steps for learning  $h$

1. Represent hypothesis class  $\mathcal{H}$
2. Define a loss objective
3. **Minimize loss to find best  $h \in \mathcal{H}$**

Many hypothesis in  $\mathcal{H}$ .

How to find the best hypothesis? – Step 3

# Optimizing Loss Functions – Step 3

- Least squares loss function

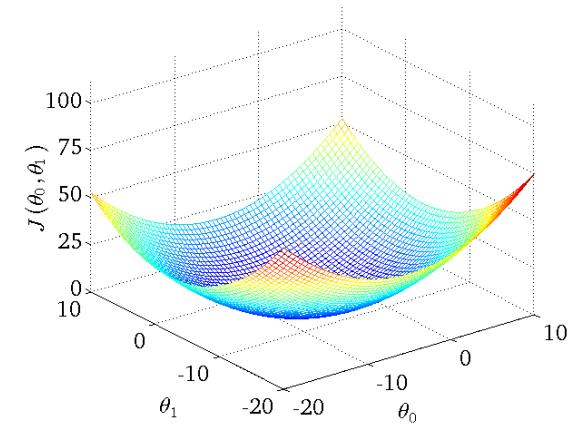
$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

- Fit by solving

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

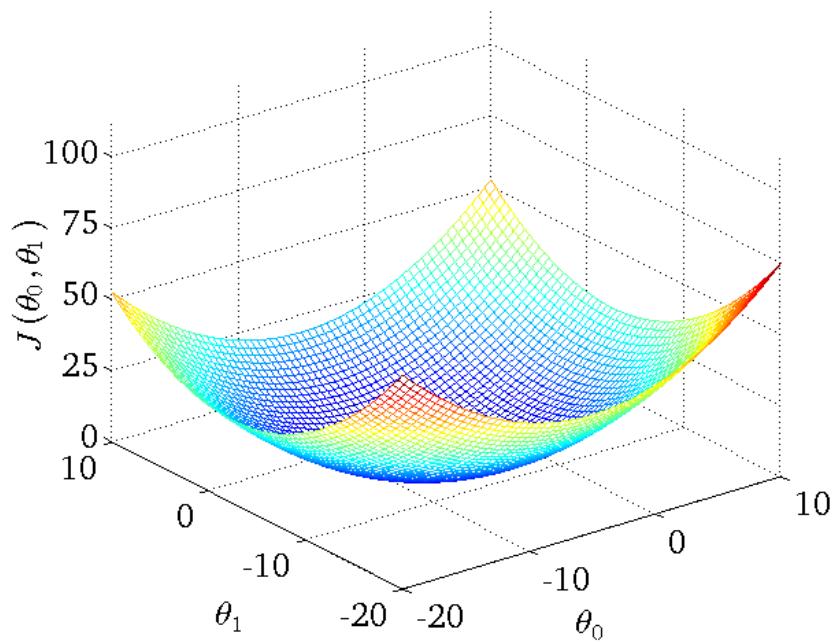
Two optimization algorithms:

- **Gradient Descent**
- **Normal Equation**



# Gradient Descent - Intuition

Think of loss at any point as its height above ground level



How will you minimize it?

# Gradient Descent - Basic Search Procedure

- Choose initial value for  $\theta$
- Until we reach a minimum:
  - Choose a new value for  $\theta$  to reduce  $J(\theta)$

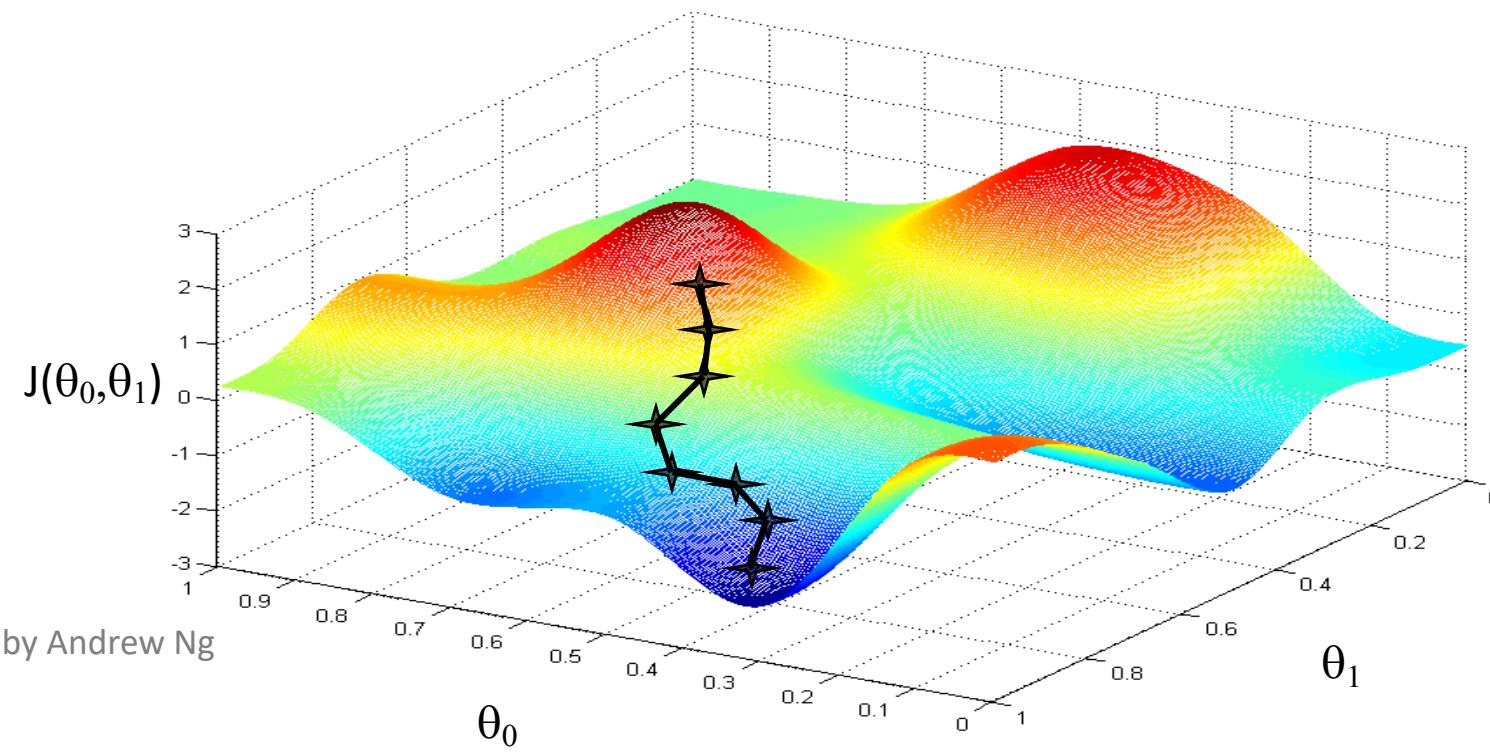
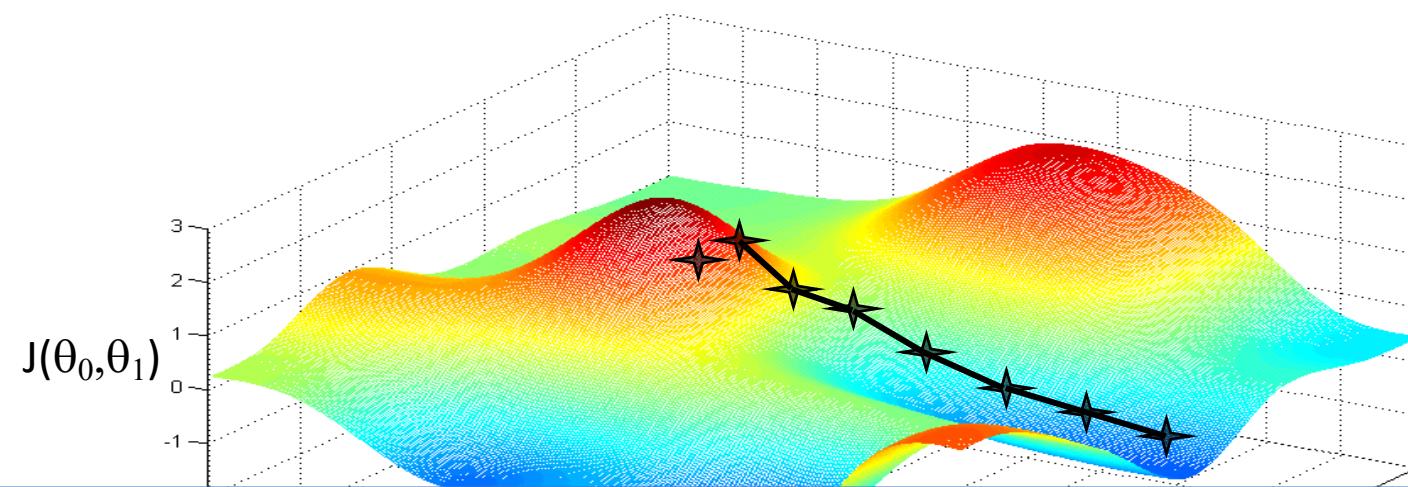


Figure by Andrew Ng

# Basic Search Procedure

- Choose initial value for  $\theta$
- Until we reach a minimum:
  - Choose a new value for  $\theta$  to reduce  $J(\theta)$



Since the least squares loss function is convex,  
we don't need to worry about local minima

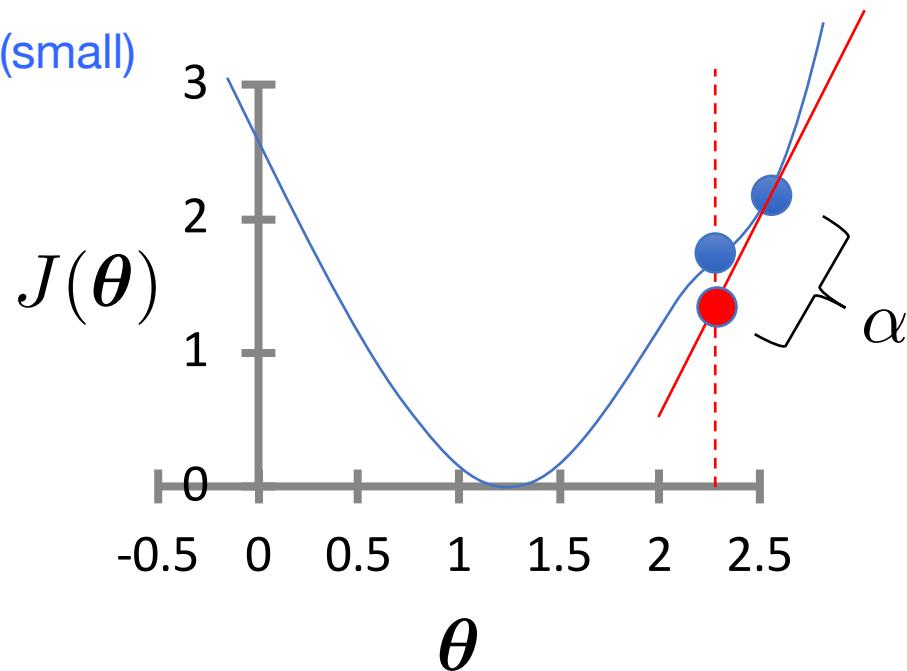
# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

simultaneous update  
for  $j = 0 \dots d$

step size or learning rate (small)  
e.g.,  $\alpha = 0.05$



# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \begin{matrix} \text{simultaneous update} \\ \text{for } j = 0 \dots d \end{matrix}$$

For Linear Regression:

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta} (x^{(i)}) - y^{(i)} \right)^2$$

# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \begin{matrix} \text{simultaneous update} \\ \text{for } j = 0 \dots d \end{matrix}$$

For Linear Regression:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta} (\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2 \end{aligned}$$

# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \begin{matrix} \text{simultaneous update} \\ \text{for } j = 0 \dots d \end{matrix}$$

For Linear Regression:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta} (\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \end{aligned}$$

# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \begin{matrix} \text{simultaneous update} \\ \text{for } j = 0 \dots d \end{matrix}$$

For Linear Regression:

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( h_{\theta} (\mathbf{x}^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)^2 \\ &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \\ &= \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)} \end{aligned}$$

# Gradient Descent

- Initialize  $\theta$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \quad \begin{matrix} \text{simultaneous update} \\ \text{for } j = 0 \dots d \end{matrix}$$

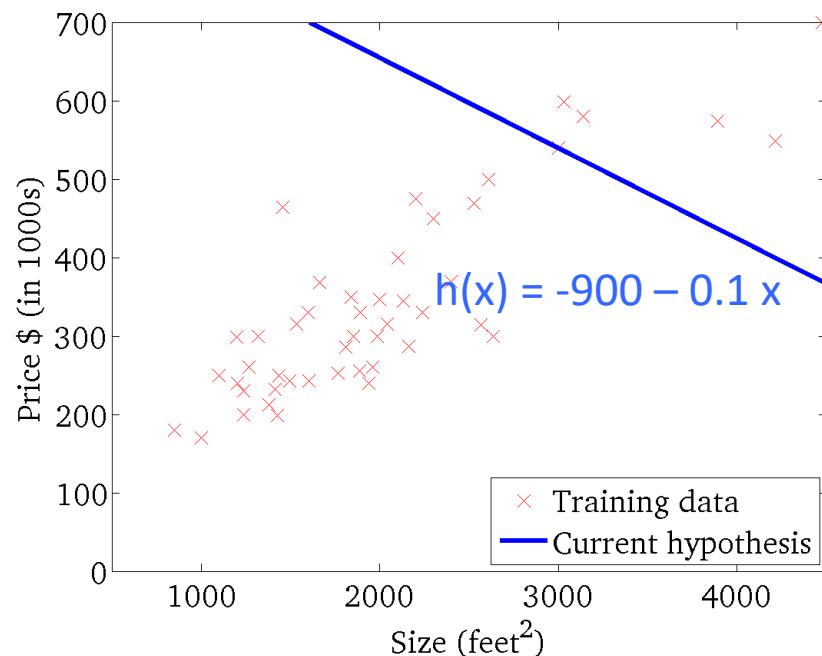
- To achieve simultaneous update
  - At the start of each GD iteration, compute  $h_{\theta}(x^{(i)})$
  - Use this stored value in the update step loop
- Assume convergence when  $\|\theta_{new} - \theta_{old}\|_2 < \epsilon$

L<sub>2</sub> norm:  $\|v\|_2 = \sqrt{\sum_i v_i^2} = \sqrt{v_1^2 + v_2^2 + \dots + v_{|v|}^2}$

# Gradient Descent

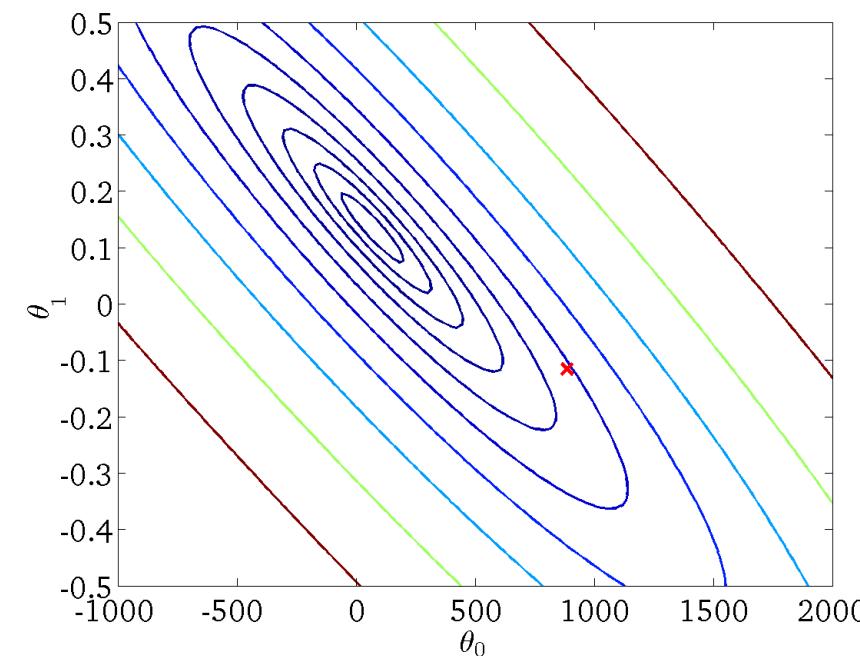
$$h_{\theta}(x)$$

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



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

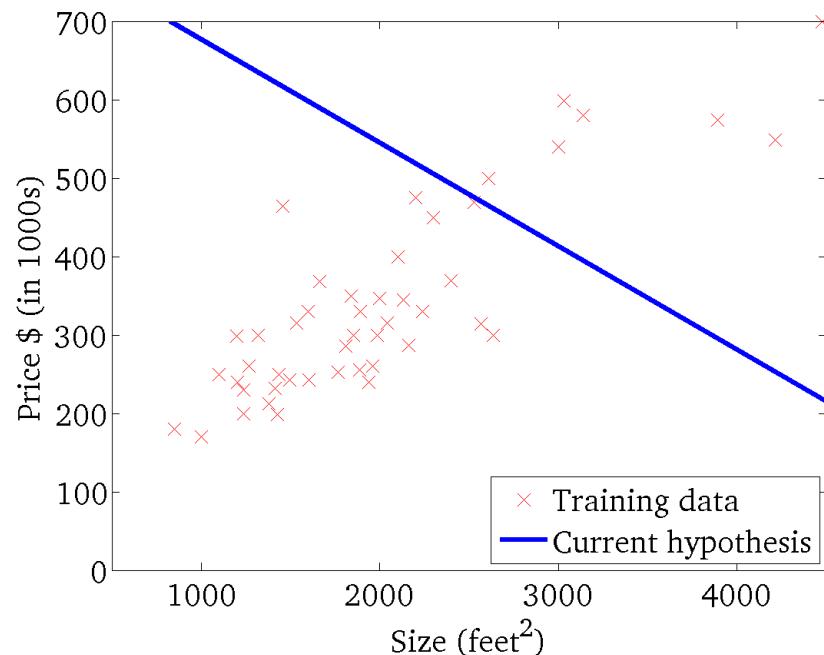
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

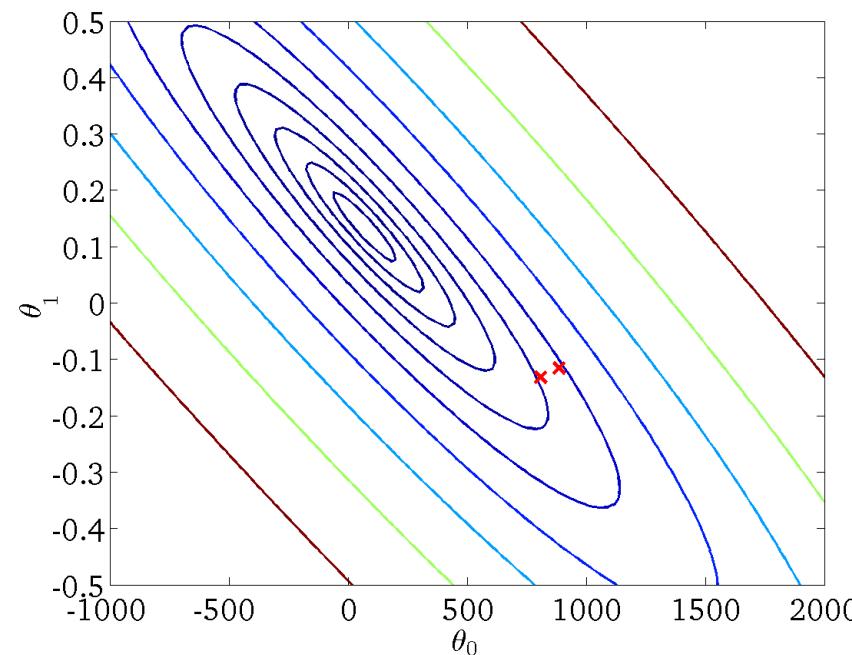
$$h_{\theta}(x)$$

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



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

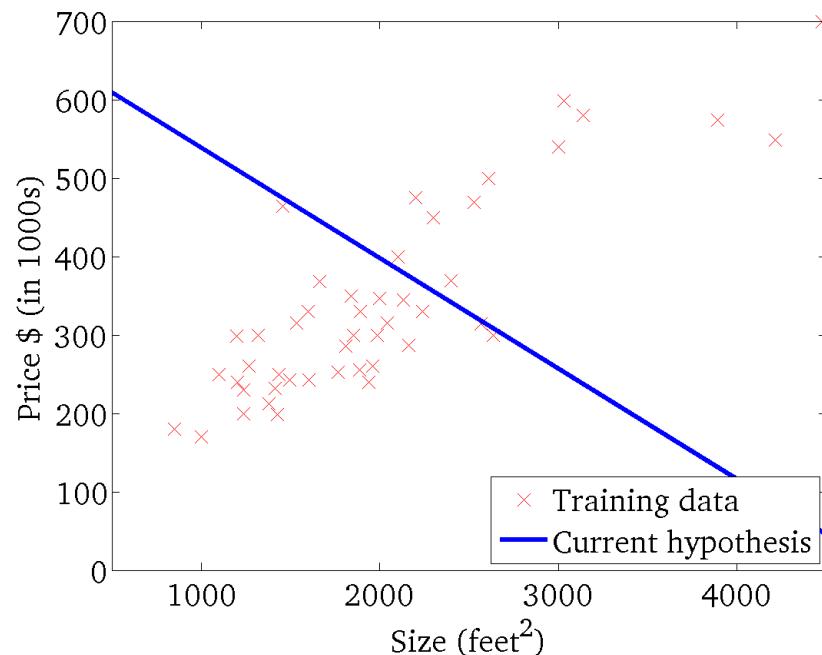
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

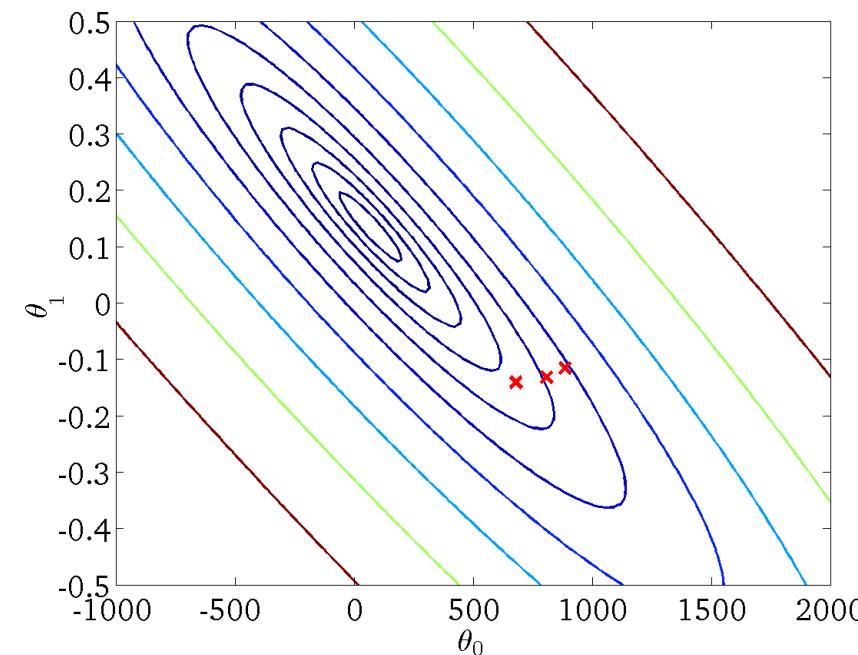
$$h_{\theta}(x)$$

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



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

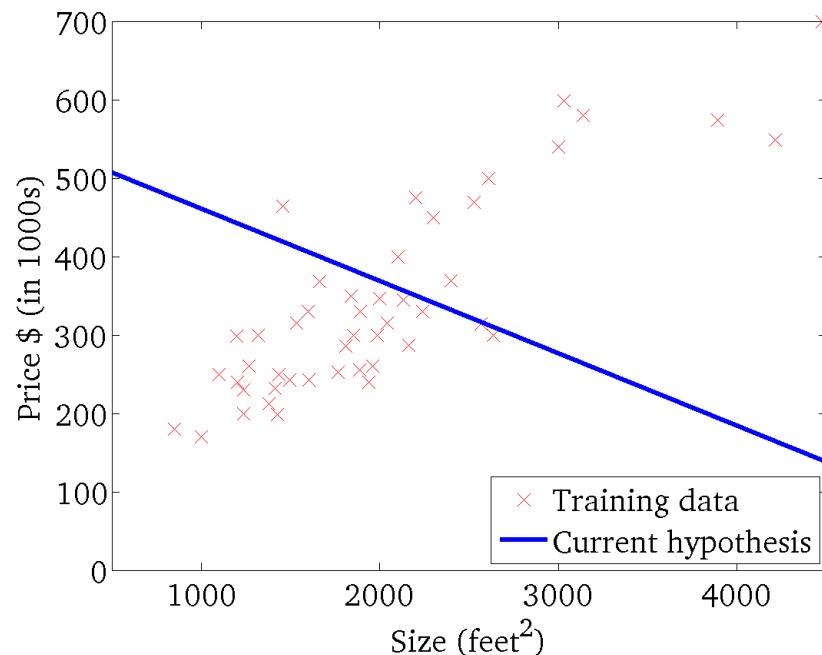
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

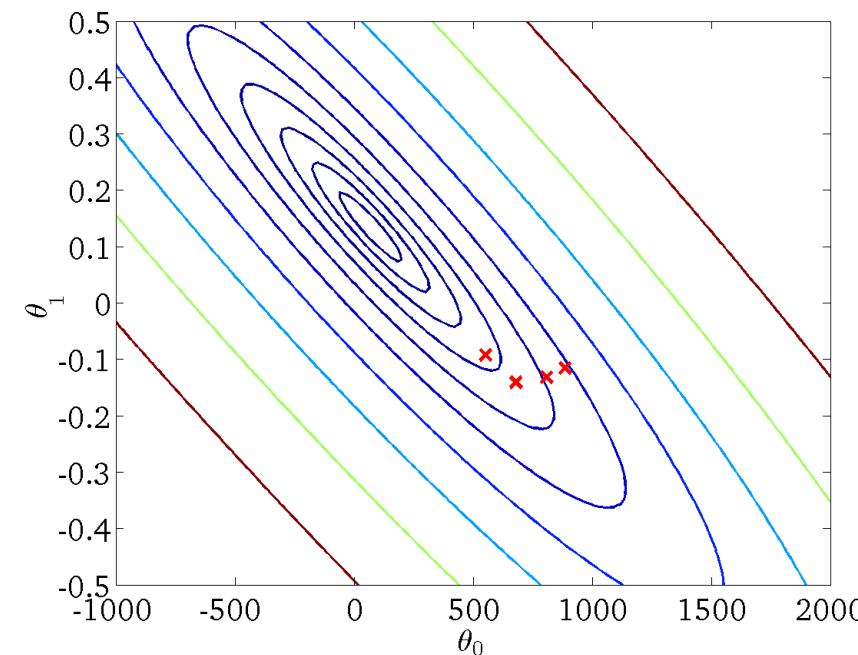
$$h_{\theta}(x)$$

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



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

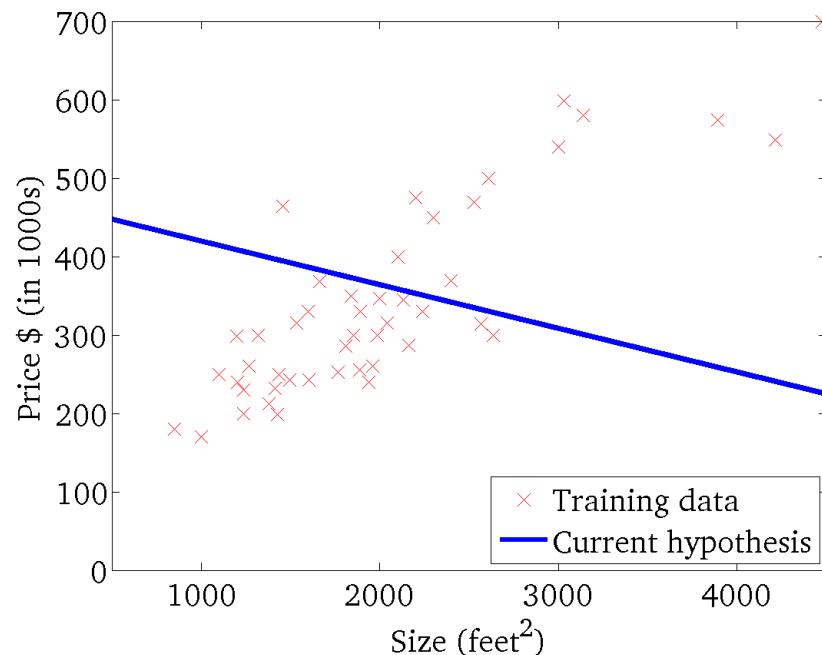
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

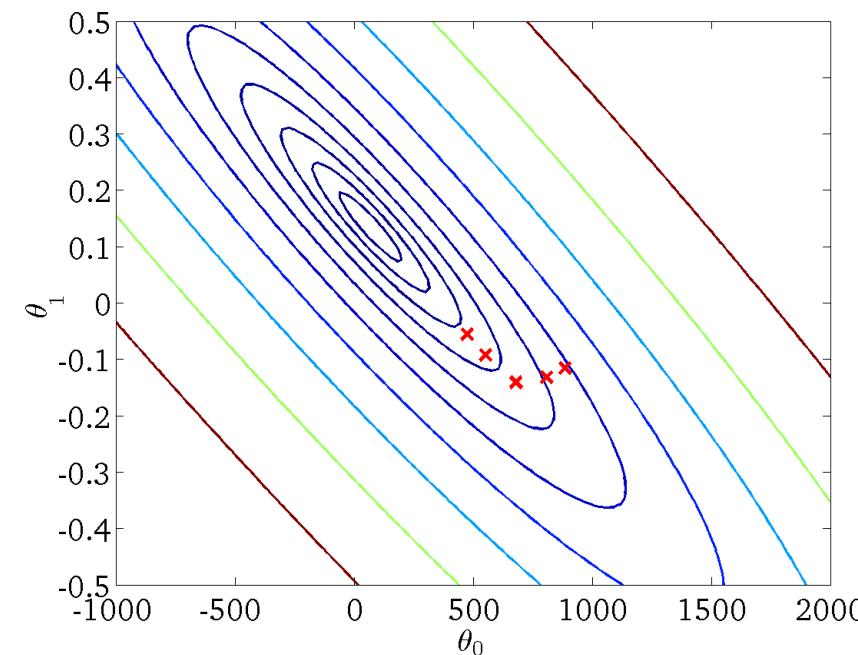
$$h_{\theta}(x)$$

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



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

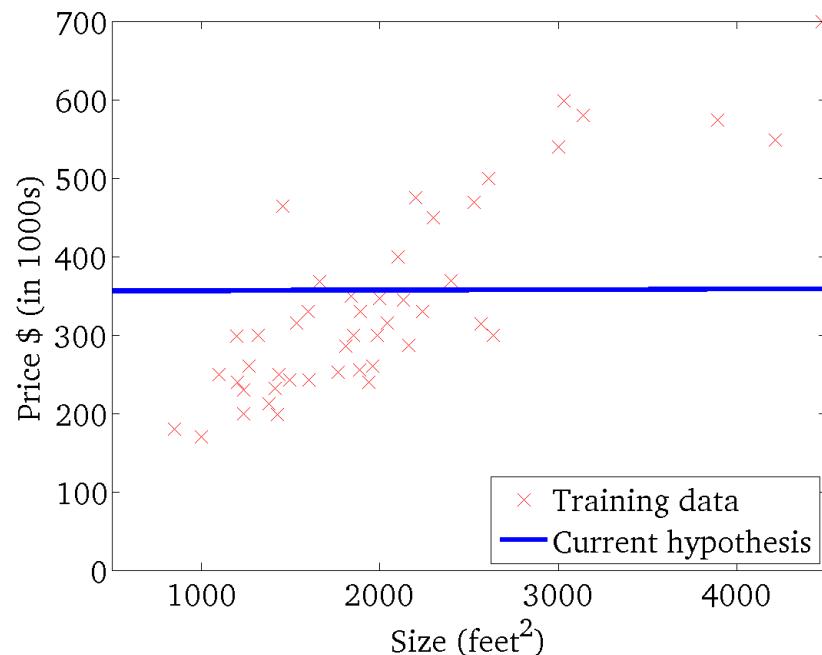
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

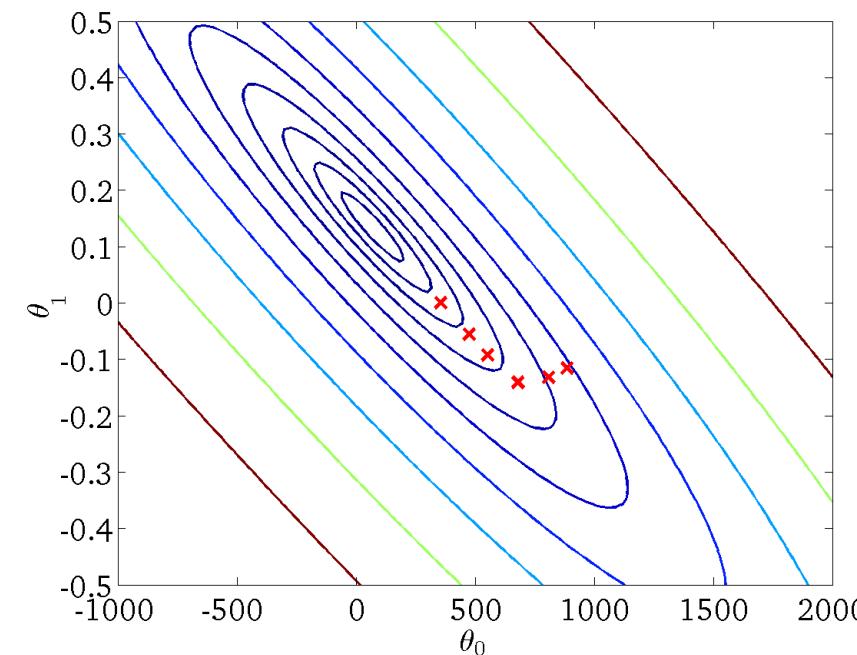
$$h_{\theta}(x)$$

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



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

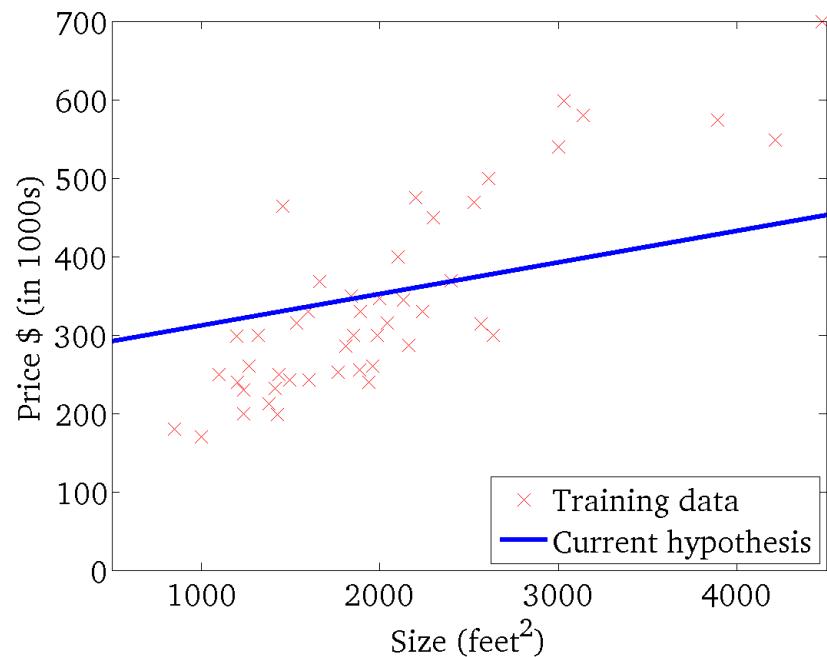
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

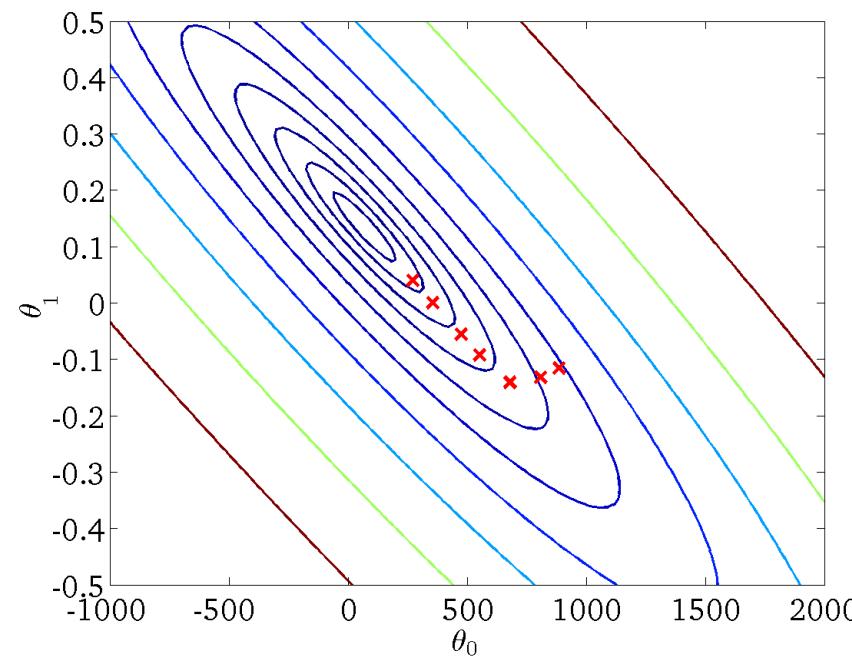
$$h_{\theta}(x)$$

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



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

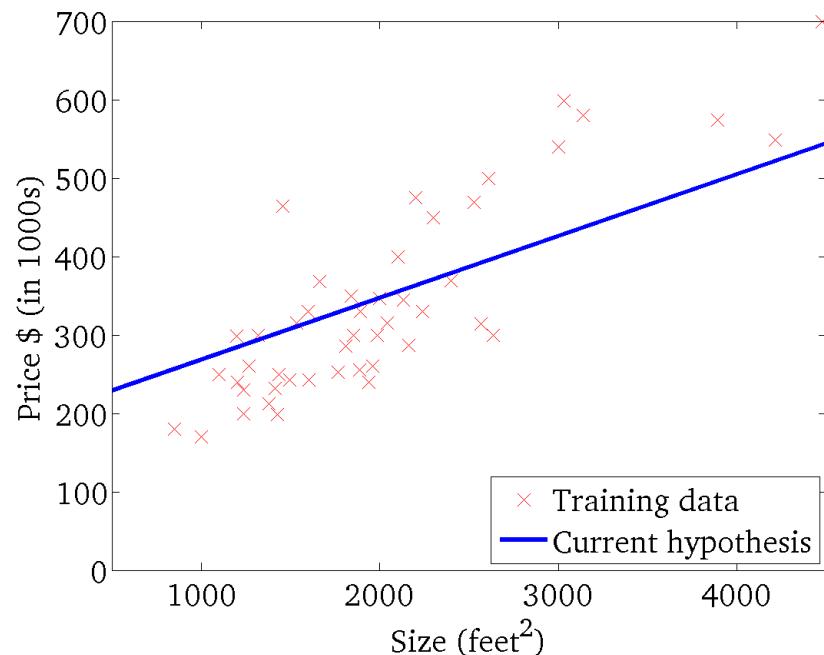
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

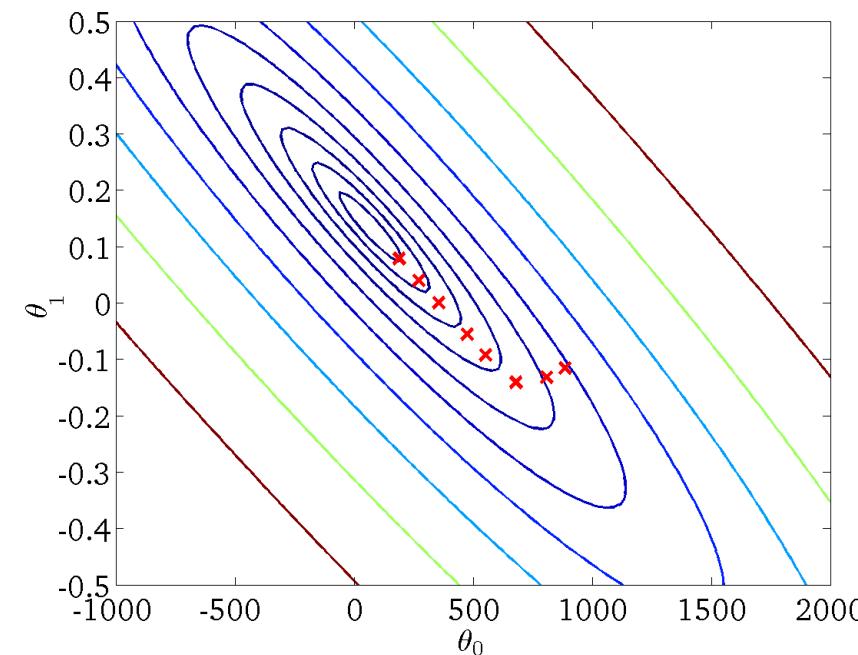
$$h_{\theta}(x)$$

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



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

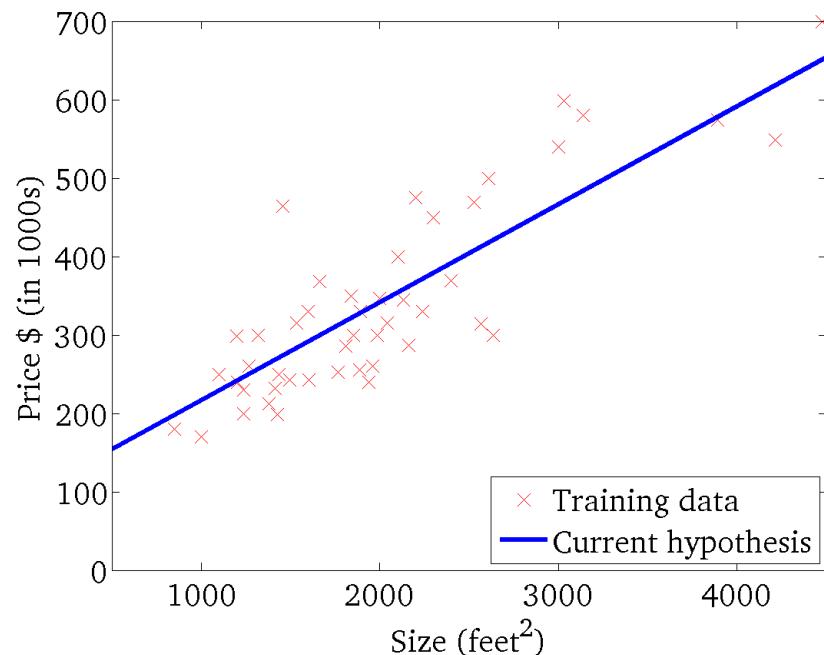
(function of the parameters  $\theta_0, \theta_1$ )



# Gradient Descent

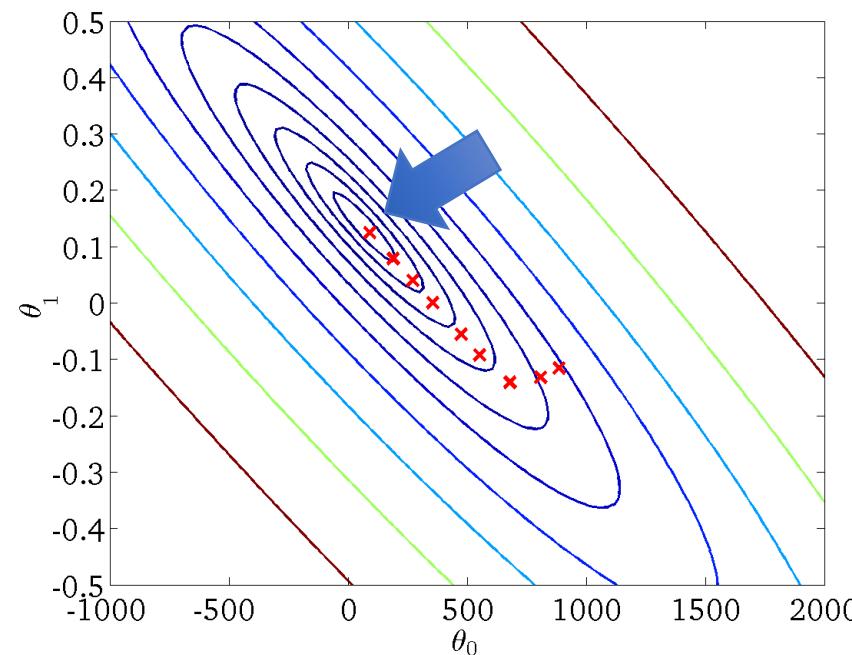
$$h_{\theta}(x)$$

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



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

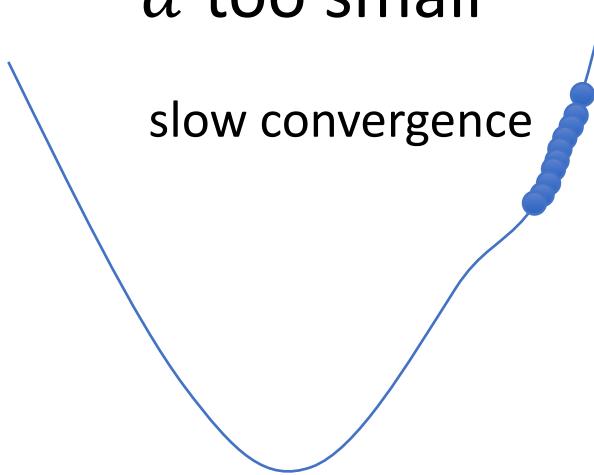
(function of the parameters  $\theta_0, \theta_1$ )



# Learning rate $\alpha$

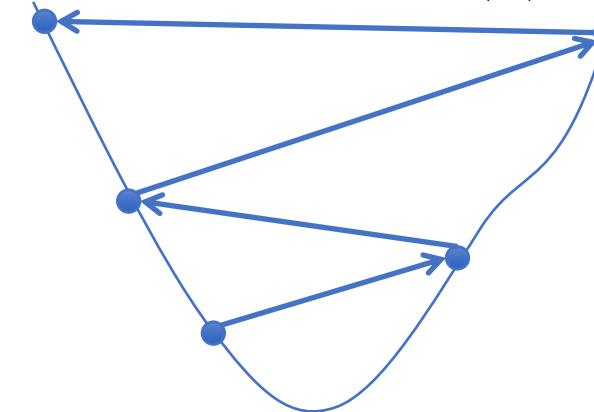
$\alpha$  too small

slow convergence



$\alpha$  too large

Increasing value for  $J(\theta)$



- May overshoot the minimum
- May fail to converge
- May even diverge

To see if gradient descent is working, print out  $J(\theta)$  each iteration

- The value should decrease at each iteration
- If it doesn't, adjust  $\alpha$

# Heuristic for setting learning rate $\alpha$

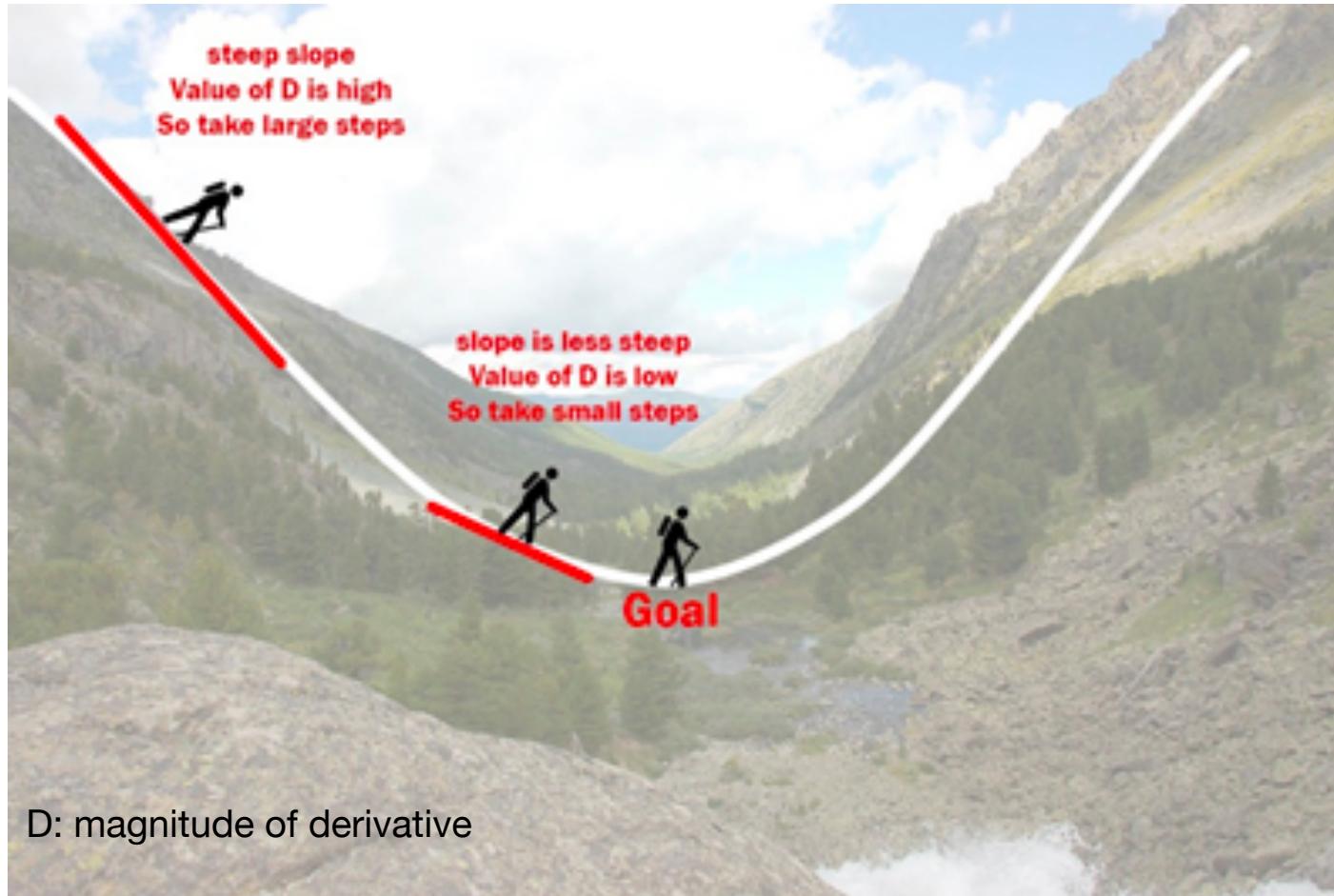


Fig src: Aarthi Kasirajan

# Scaling to large datasets

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

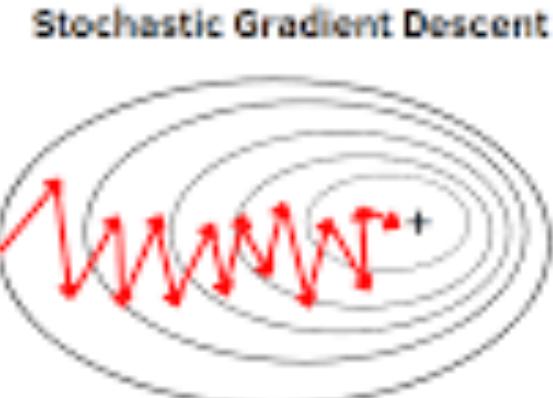
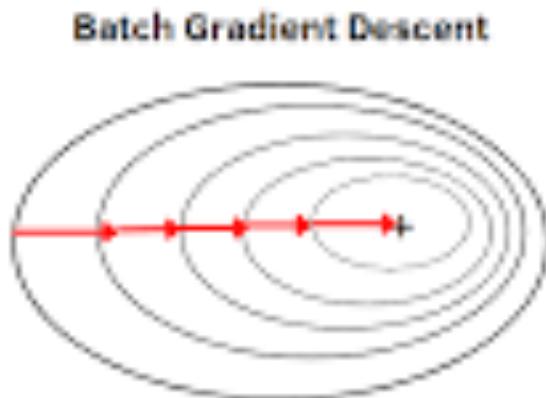
where  $\ell(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$  is the mean-squared error of  $h_{\boldsymbol{\theta}}(\mathbf{x})$  for  $(\mathbf{x}^{(i)}, y^{(i)})$

- Modern datasets for ML are very large
  - E.g., vision and language datasets have billions of training instances
- Computing average gradient over  $n$  training instances requires  $O(n)$  time
  - Very slow for an iterative algorithm when  $n$  is large!
  - Also known as **Full Batch Gradient Descent (GD)**
- **Stochastic gradient descent (SGD):** compute gradient for every example

Iterate over  $i = 1, 2, \dots, n$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} \ell(\mathbf{x}^{(i)}, y^{(i)}, \boldsymbol{\theta})$$

# GD vs SGD



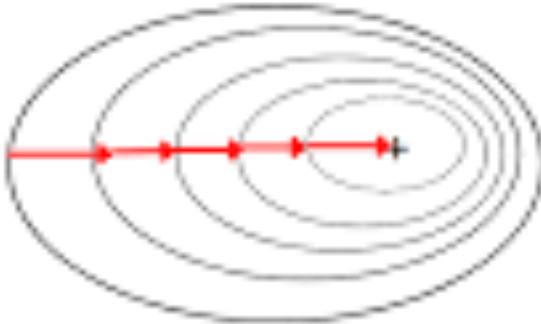
- Batch gradient descent
- Stochastic gradient descent

$$\theta \leftarrow \theta - \alpha \frac{1}{n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

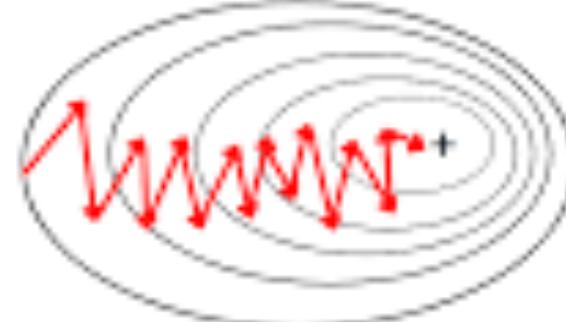
$$\theta \leftarrow \theta - \alpha (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}$$

# GD vs SGD

Batch Gradient Descent



Stochastic Gradient Descent



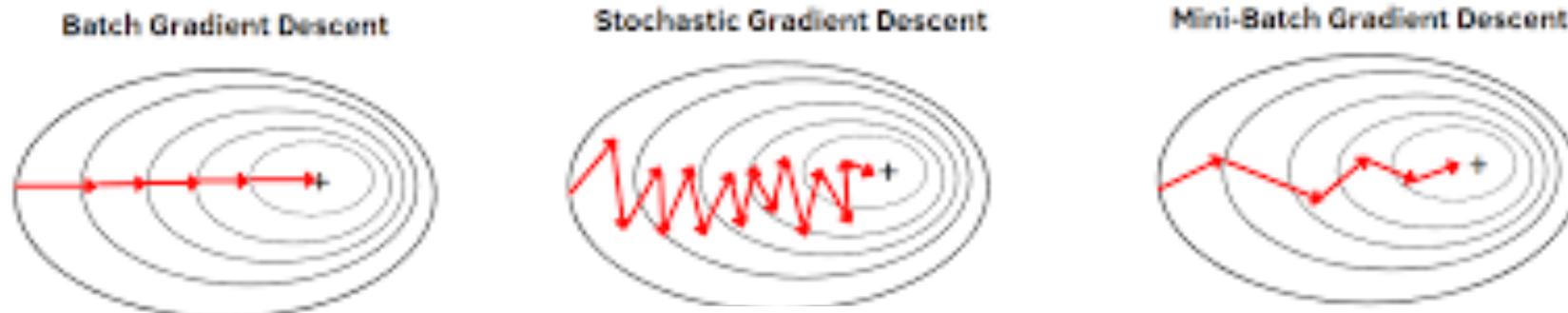
- Advantages of SGD:  
Scalability to large datasets
  - Memory efficient: single example is processed at a time
  - Computationally cheap: single example needed to compute gradient
  - Implicit regularization (advanced ML course)
- Disadvantages of SGD
  - High noise in the gradient. Can take many updates to converge
  - Cannot exploit modern hardware optimized for matrix operations

# Mini-batch Gradient Descent

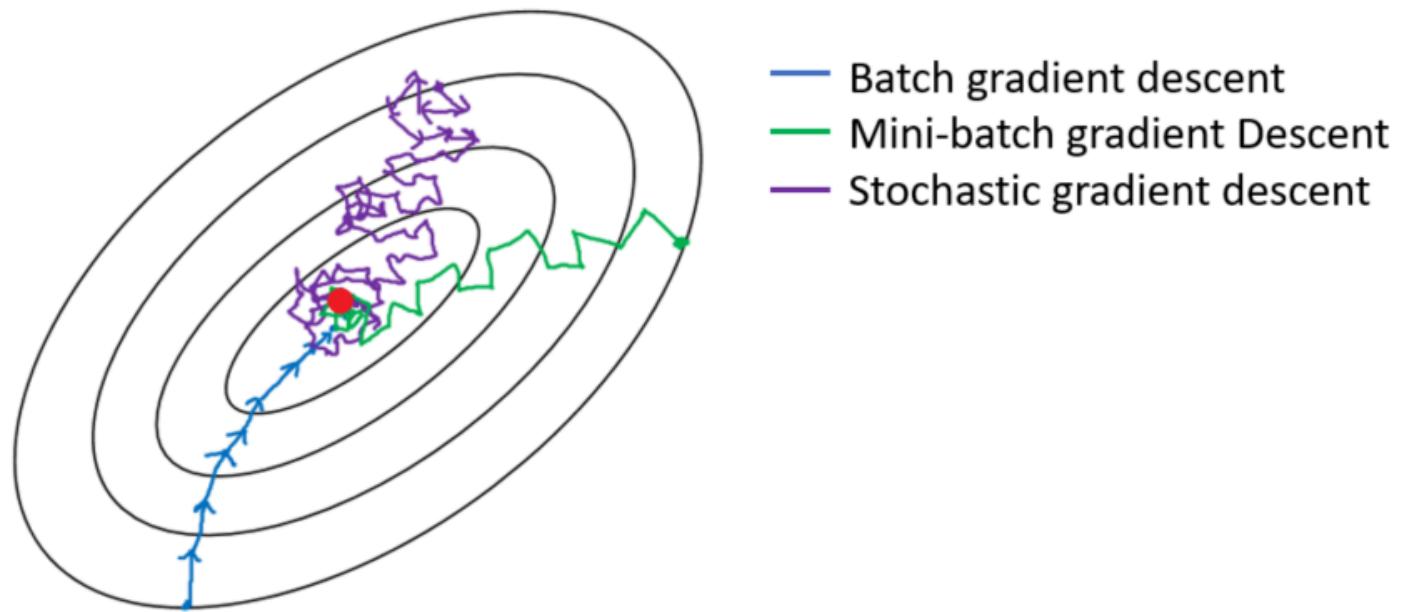
- **Mini-batch GD (MBGD):** Sample a batch of  $B$  points  $D_B$  at random from the full dataset  $D$  without replacement

$$\theta \leftarrow \theta - \alpha \frac{1}{B} \sum_{(x^{(i)}, y^{(i)}) \in D_B} \nabla_{\theta} \ell(x^{(i)}, y^{(i)}, \theta)$$

- Compromise between (full batch) GD and SGD
  - When  $B = 1$ , MBGD is same as SGD
  - When  $B = n$ , MBGD is same as GD



# GD vs SGD vs MBGD



# Optimizing Loss Functions – Step 3

- Least squares loss function

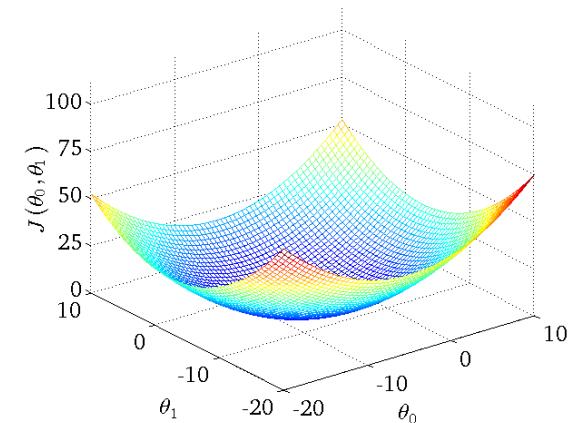
$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2$$

- Fit by solving

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

Two optimization algorithms:

- **Gradient Descent**
- **Normal Equation** – next!



# Vectorization

- Benefits of vectorization
  - More compact equations
  - Faster code (using optimized matrix libraries)
- Consider our model:

$$h(\mathbf{x}) = \sum_{j=0}^d \theta_j x_j$$

- Let

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = [ 1 \quad x_1 \quad \dots \quad x_d ]$$

- Can write the model in vectorized form as

$$h(\mathbf{x}) = \boldsymbol{\theta}^\top \mathbf{x}$$

# Vectorization

- Consider our model for  $n$  instances:

$$h(\mathbf{x}^{(i)}) = \sum_{j=0}^d \theta_j x_j^{(i)}$$

- Let

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad \mathbb{R}^{n \times (d+1)}$$

- Can write the model in vectorized form as

$$h_{\theta}(\mathbf{x}) = \mathbf{X}\theta$$

# Vectorization

- For the linear regression cost function:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

$$= \frac{1}{2n} \sum_{i=1}^n \left( \boldsymbol{\theta}^\top \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{2n} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$$

Let:

$$\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

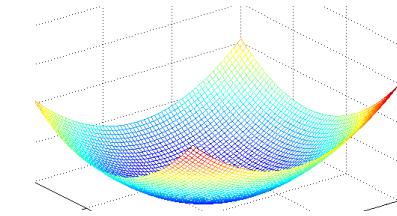
The diagram illustrates the dimensions of the terms in the vectorized cost function. It shows the expression  $\frac{1}{2n} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$ . Blue arrows point from each term to its corresponding dimension:

- A blue arrow points from  $\mathbf{X}\boldsymbol{\theta}$  to  $\mathbb{R}^{n \times (d+1)}$ .
- A blue arrow points from  $\mathbf{y}$  to  $\mathbb{R}^{(d+1) \times 1}$ .
- A blue bracket under  $\mathbf{X}\boldsymbol{\theta} - \mathbf{y}$  indicates a  $\mathbb{R}^{1 \times n}$  vector.
- A blue bracket under  $\mathbf{X}\boldsymbol{\theta} - \mathbf{y}$  indicates a  $\mathbb{R}^{n \times 1}$  vector.

# Normal Equations - Closed Form Solution

- Instead of using GD, solve for optimal  $\theta$  analytically

$$\frac{\partial}{\partial \theta} J(\theta) = 0$$



- Derivation:

$$\begin{aligned}\mathcal{J}(\theta) &= \frac{1}{2n} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y}) \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - \boxed{\mathbf{y}^\top \mathbf{X} \theta} - \boxed{\theta^\top \mathbf{X}^\top \mathbf{y}} + \mathbf{y}^\top \mathbf{y} \\ &\propto \theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}\end{aligned}$$

1 x 1

Take derivative and set equal to 0, then solve for  $\theta$ :

$$\begin{aligned}\frac{\partial}{\partial \theta} (\theta^\top \mathbf{X}^\top \mathbf{X} \theta - 2\theta^\top \mathbf{X}^\top \mathbf{y} + \cancel{\mathbf{y}^\top \mathbf{y}}) &= 0 \\ (\mathbf{X}^\top \mathbf{X})\theta - \mathbf{X}^\top \mathbf{y} &= 0 \\ (\mathbf{X}^\top \mathbf{X})\theta &= \mathbf{X}^\top \mathbf{y}\end{aligned}$$

Closed Form Solution:  $\theta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

# Closed Form Solution

- Can obtain  $\theta$  by simply plugging  $X$  and  $y$  into

$$\theta = (X^T X)^{-1} X^T y$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- If  $X^T X$  is not invertible (i.e., singular), may need to:
  - Use pseudo-inverse instead of the inverse
    - In python, `numpy.linalg.pinv(a)`
  - Remove redundant (not linearly independent) features
  - Remove extra features to ensure that  $d \leq n$

# Gradient Descent v.s. Normal Equations

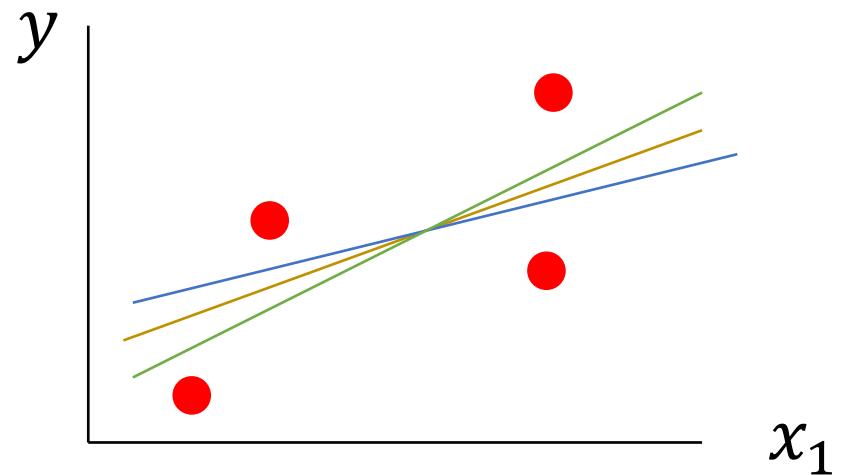
## Gradient Descent

- Requires multiple iterations
- Need to choose learning rate  $\alpha$
- Works well when  $n$  is large
- Can support incremental learning

## Normal Equations

- Non-iterative
- No need for  $\alpha$
- Slow if  $d$  is large
  - Computing  $(\mathbf{X}^\top \mathbf{X})^{-1}$  is roughly  $O(d^3)$

# Linear Regression



Steps for learning  $h$

1. Represent hypothesis class  $\mathcal{H}$
2. Define a loss objective
3. Minimize loss to find best  $h \in \mathcal{H}$

# Summary so far

## 3 Steps:

- Represent hypothesis class

**Linear functions:**  $h_{\theta}(x) = \theta^T x$  (linear in both  $\theta$  and  $x$ )

- Define a loss function

**Squared error loss:**  $J(\theta) = \frac{1}{2n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2$

- Optimize loss to find best hypothesis

**Gradient descent:**  $\theta \leftarrow \theta - \alpha \frac{1}{n} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$

**Normal equations:**  $\hat{\theta} = (X^T X)^{-1} X^T y$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(i)} & \dots & x_d^{(i)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$