

Linear Regression

Fundamentals of Data Science
20, 23, 27, 30 October, 3, 6 November 2023
Prof. Fabio Galasso



Outline

- Linear regression
 - One or multiple variables
 - Cost function
 - Gradient descent, incl. batch/stochastic GD
 - Normal equation
 - MSE and Correlation
 - Locally-Weighted Regression
 - Probabilistic Interpretation of Least Squares
-

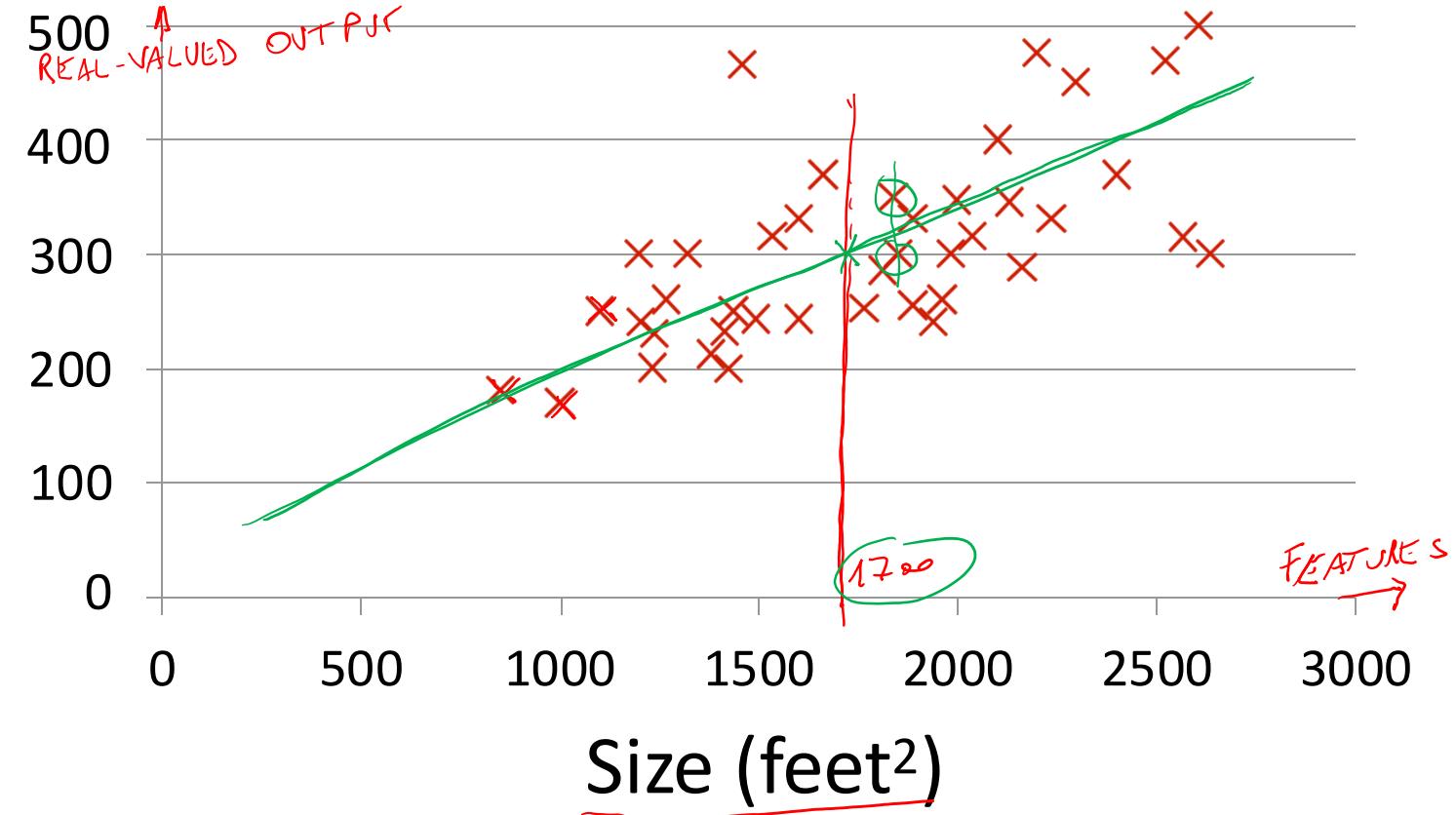
Outline

- Linear regression
 - One or multiple variables
 - Cost function
 - Gradient descent, incl. batch/stochastic GD
 - Normal equation
 - MSE and Correlation
 - Locally-Weighted Regression
 - Probabilistic Interpretation of Least Squares
-

Linear Regression

Housing Prices (Portland, OR)

Price
(in 1000s
of dollars)



Supervised Learning

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

Regression Problem

Predict real-valued output

CLASSIFICATION : DISCRETE-VALUED OUTPUT

Training set of housing prices (Portland, OR)

	Size in feet ² (x)	Price (\$) in 1000's (y)
1)	2104	460
2)	1416	232
3)	1534	315
4)	852	178

Notation:

m = Number of training examples

x 's = “input” variable / features

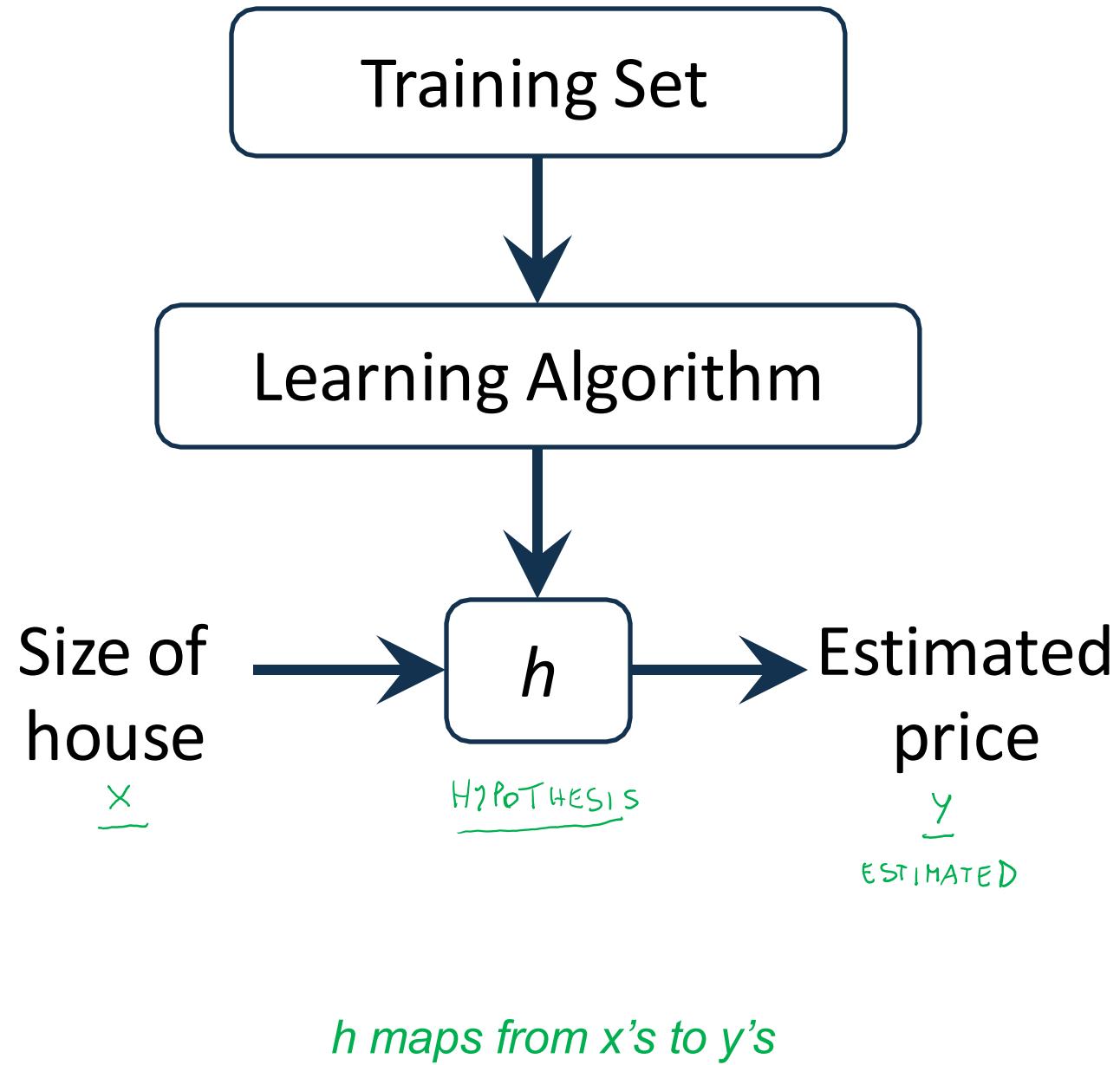
y 's = “output” variable / “target” variable

$$x^{(2)} = 1416$$

$$y^{(2)} = 232$$

(x, y) TRAINING EXAMPLE

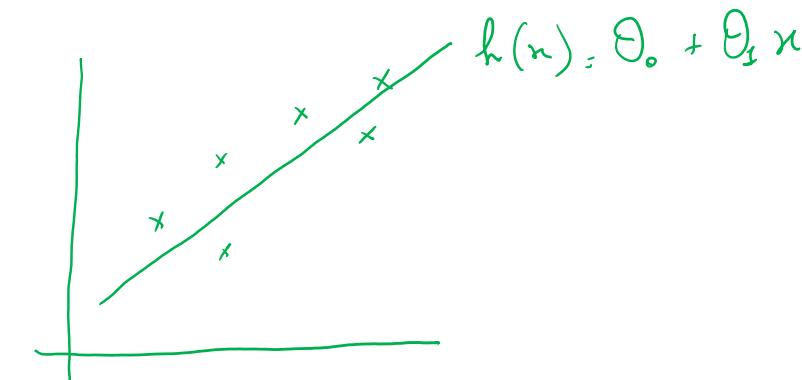
$(x^{(i)}, y^{(i)})$ i -th TRAINING EXAMPLE



How do we represent h ?

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

SHORTHAND $h(x)$



Linear regression with one variable. (x)
Univariate linear regression.

Multiple features (variables)

Size (feet ²)	Price (\$1000)
x	y
2104	460
1416	232
1534	315
852	178
...	...

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

Multiple features (variables)

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
<u>1416</u>	<u>3</u>	<u>2</u>	<u>40</u>	<u>232</u>
1534	3	2	30	315
852	2	1	36	178
...

Notation:

→ n = number of features

→ $x^{(i)}$ = input (features) of i^{th} training example.

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

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$$x_3^{(2)} = 2$$

m TRAINING EXAMPLES

Hypothesis:

With one variable: $h_{\theta}(x) = \theta_0 + \theta_1 x$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$h_{\theta}(w) = 100 + 0.1 w_1 + 0.6 w_2 + 0.07 w_3 + 10 w_4$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

For convenience of notation, define $\underline{x_0} = 1$. $\underline{x_0}^{(i)} = 1$

$$\underline{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \underline{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix} \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}}_{\underline{x}}$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

$$= \underline{\theta}^T \underline{x}$$

Multivariate linear regression

Linear Regression: Cost Function

Training Set

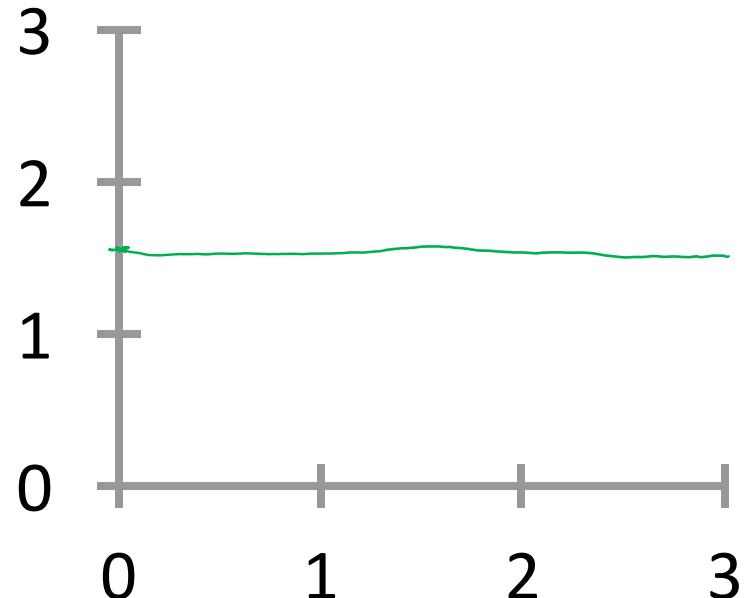
	Size in feet ² (x)	Price (\$) in 1000's (y)
	2104	460
	1416	232
	1534	315
	852	178

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

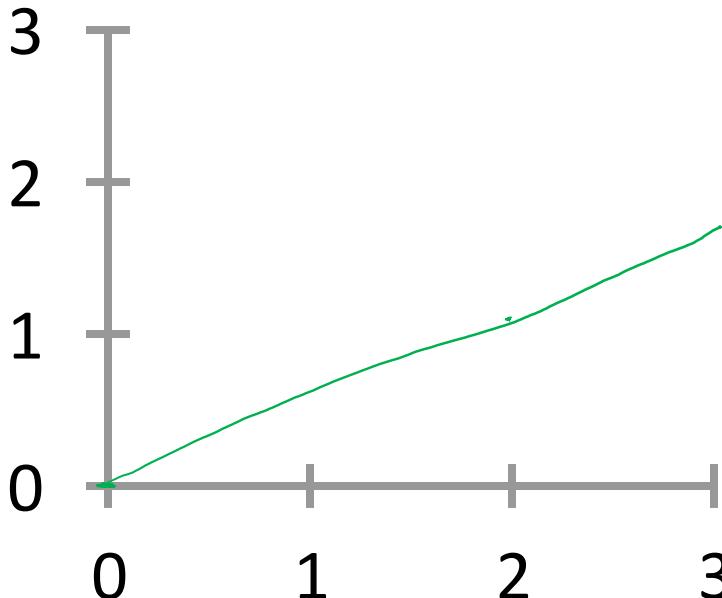
θ_i 's: Parameters

How to choose θ_i 's ?

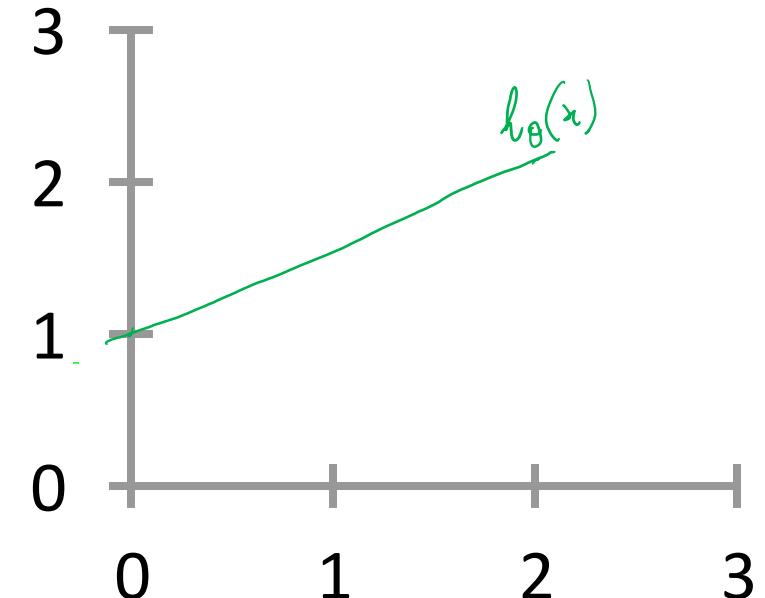
$$h_{\theta}(x) = \underline{\theta_0} + \theta_1 x$$



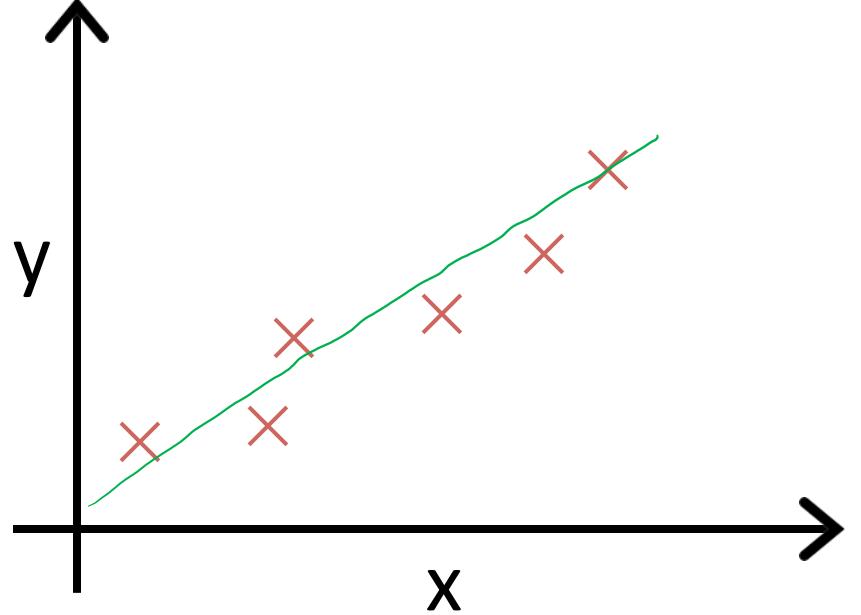
$$\begin{array}{l} \rightarrow \theta_0 = 1.5 \\ \theta_1 = 0 \end{array}$$



$$\left\{ \begin{array}{l} \theta_0 = 0 \\ \theta_1 = 0.5 \end{array} \right.$$



$$\begin{array}{l} \theta_0 = 1 \\ \theta_1 = 0.5 \end{array}$$



Idea: Choose θ_0, θ_1 so that
 $h_\theta(x)$ is close to y for our
 training examples (x, y)

minimize θ_0, θ_1
$$\frac{1}{2m} \sum_{i=1}^m \left(h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$

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

minimize θ_0, θ_1 $J(\theta_0, \theta_1)$

COST FUNCTION

Simplified

Hypothesis:

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

Parameters:

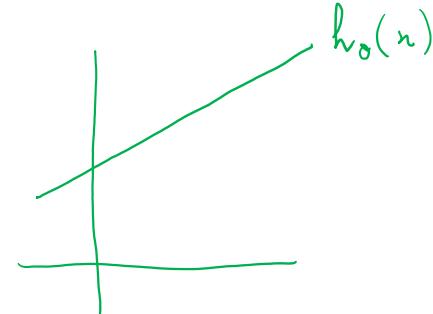
$$\theta_0, \theta_1$$

Cost Function:

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

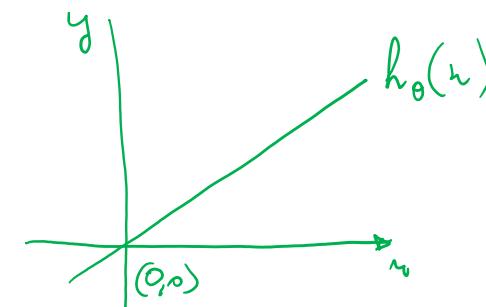
TRAINING EXAMPLES

Goal: minimize $J(\theta_0, \theta_1)$



$$h_{\theta}(x) = \theta_1 x$$

$$\theta_1$$



$$J(\underline{\theta}_1) = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

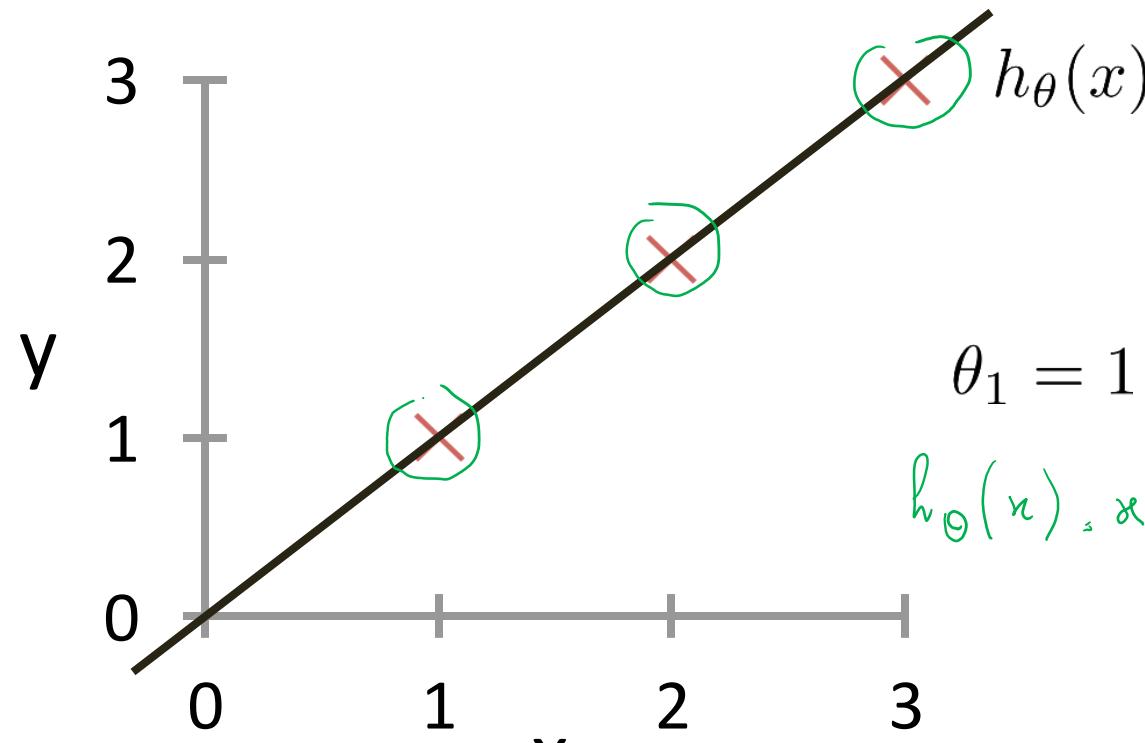
minimize $\underline{\theta}_1$ $J(\theta_1)$

$$\partial_{\theta_1} J$$

$$\partial_{\theta_1} x^{(i)}$$

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

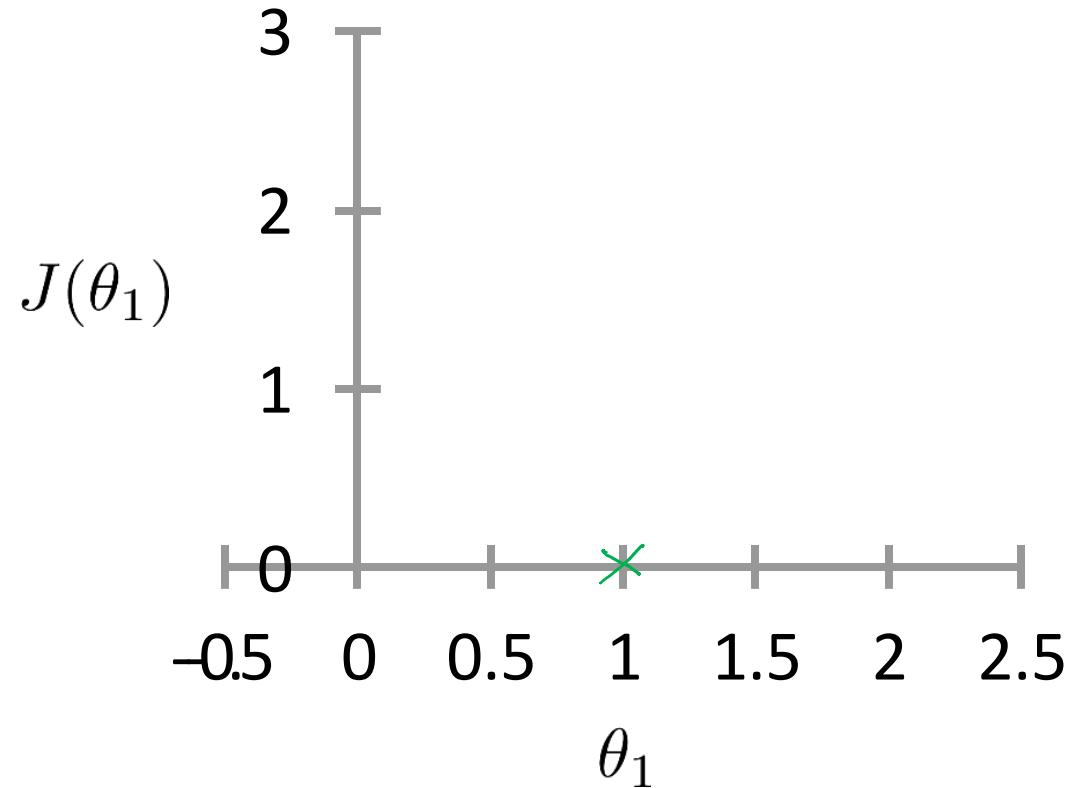
(for fixed θ_1 , this is a function of x)



$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{1}{2 \times 3} \left(0^2 + 0^2 + 0^2 \right) = 0 \end{aligned}$$

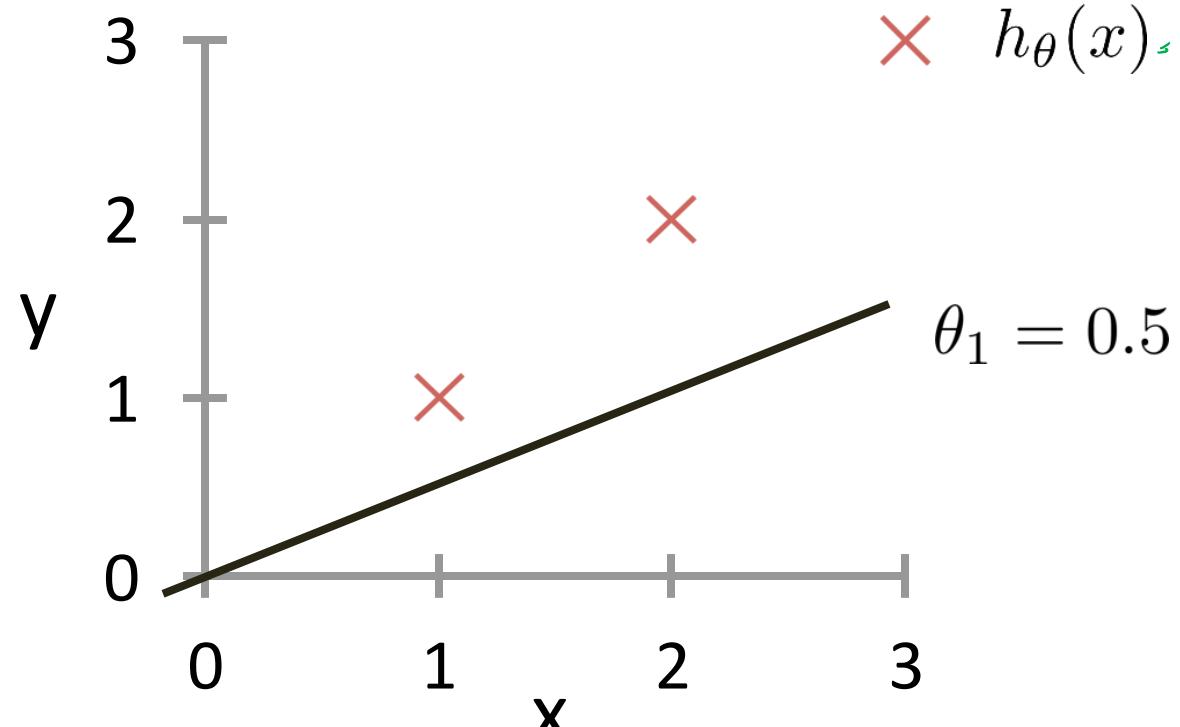
$$J(\theta_1)$$

(function of the parameter θ_1)



$h_{\theta}(x)$

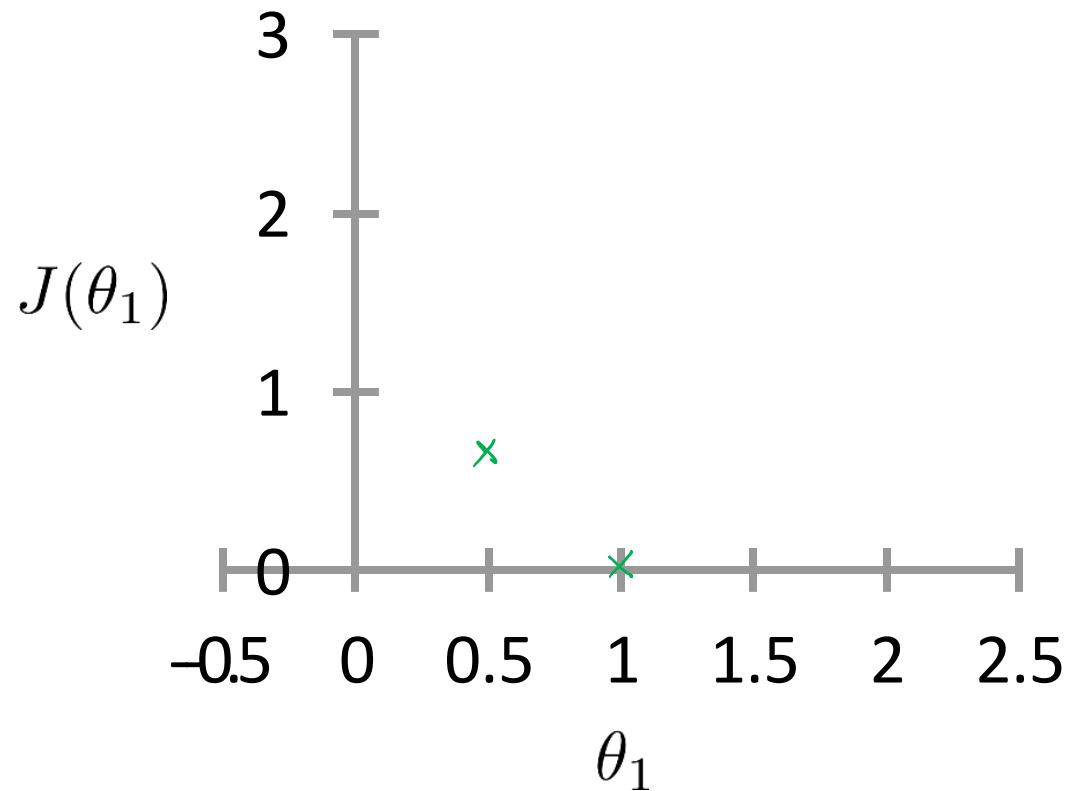
(for fixed θ_1 , this is a function of x)



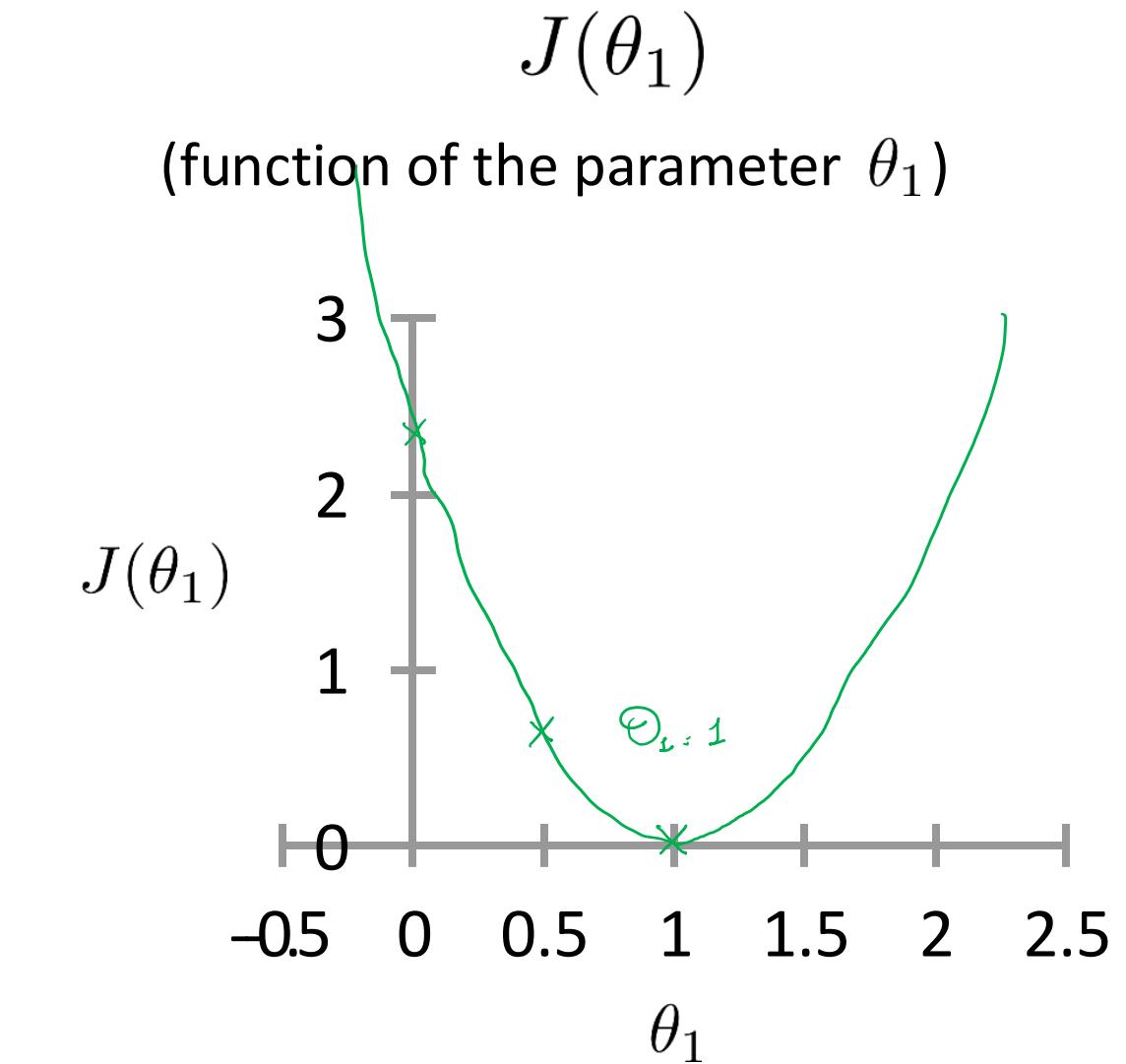
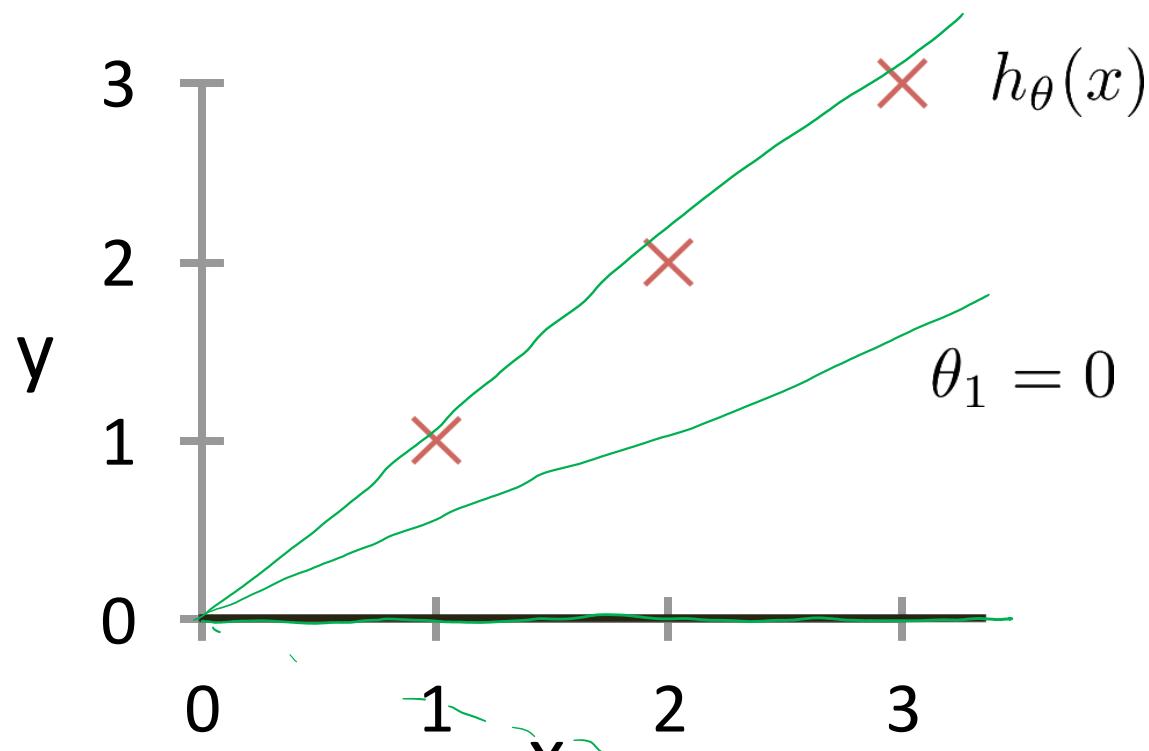
$$\begin{aligned} J(0.5) &= \frac{1}{2 \times 3} \left[(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2 \right] \\ &= \frac{1}{2 \times 3} \left(\frac{1}{4} + 1 + \frac{9}{4} \right) \simeq 0.58 \end{aligned}$$

$J(\theta_1)$

(function of the parameter θ_1)



$h_\theta(x)$
 (for fixed θ_1 , this is a function of x)



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

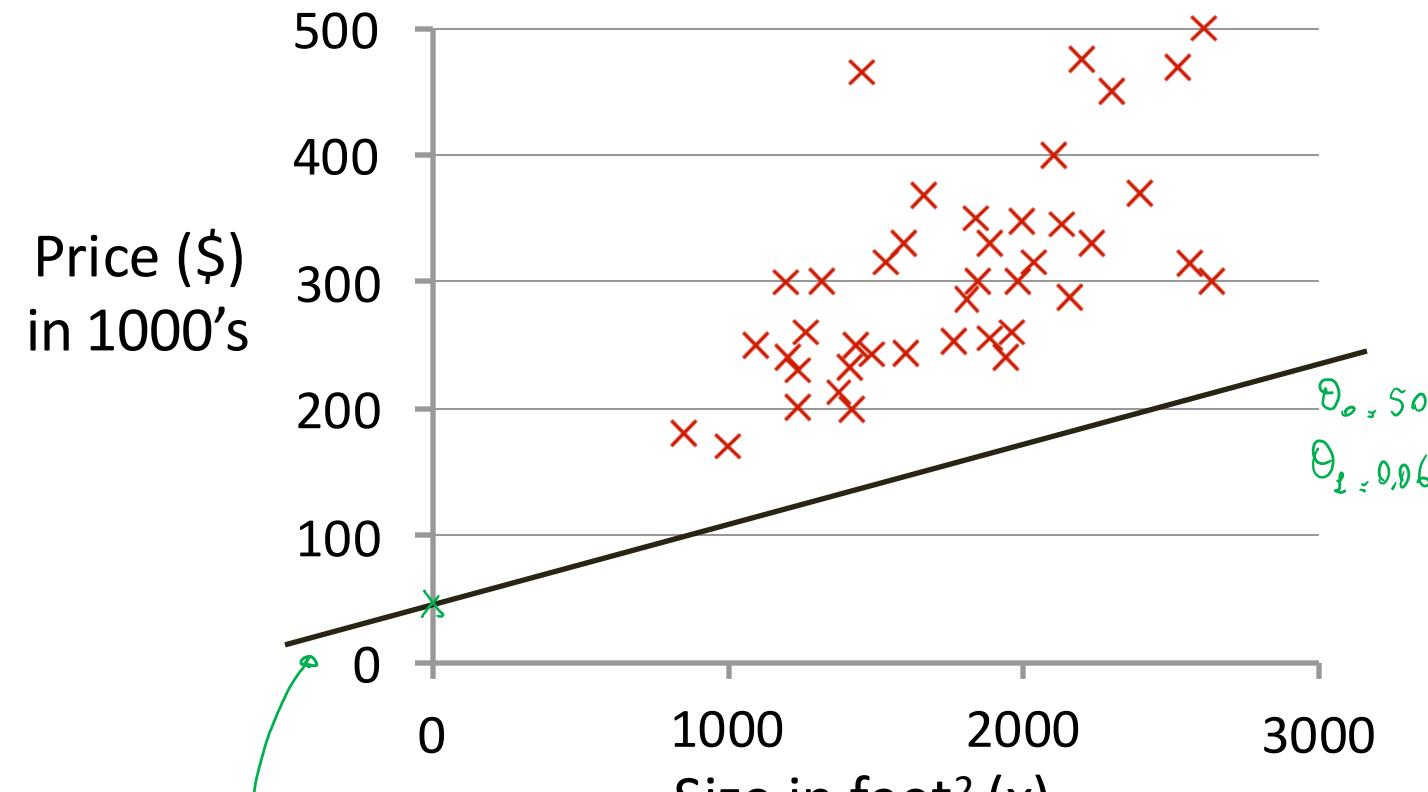
Parameters: θ_0, θ_1

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

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

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

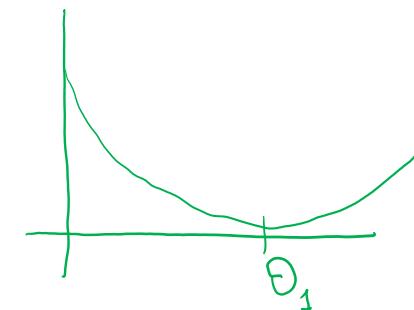
(for fixed θ_0, θ_1 , this is a function of x)



$$h_{\theta}(x) = 50 + 0.06x$$

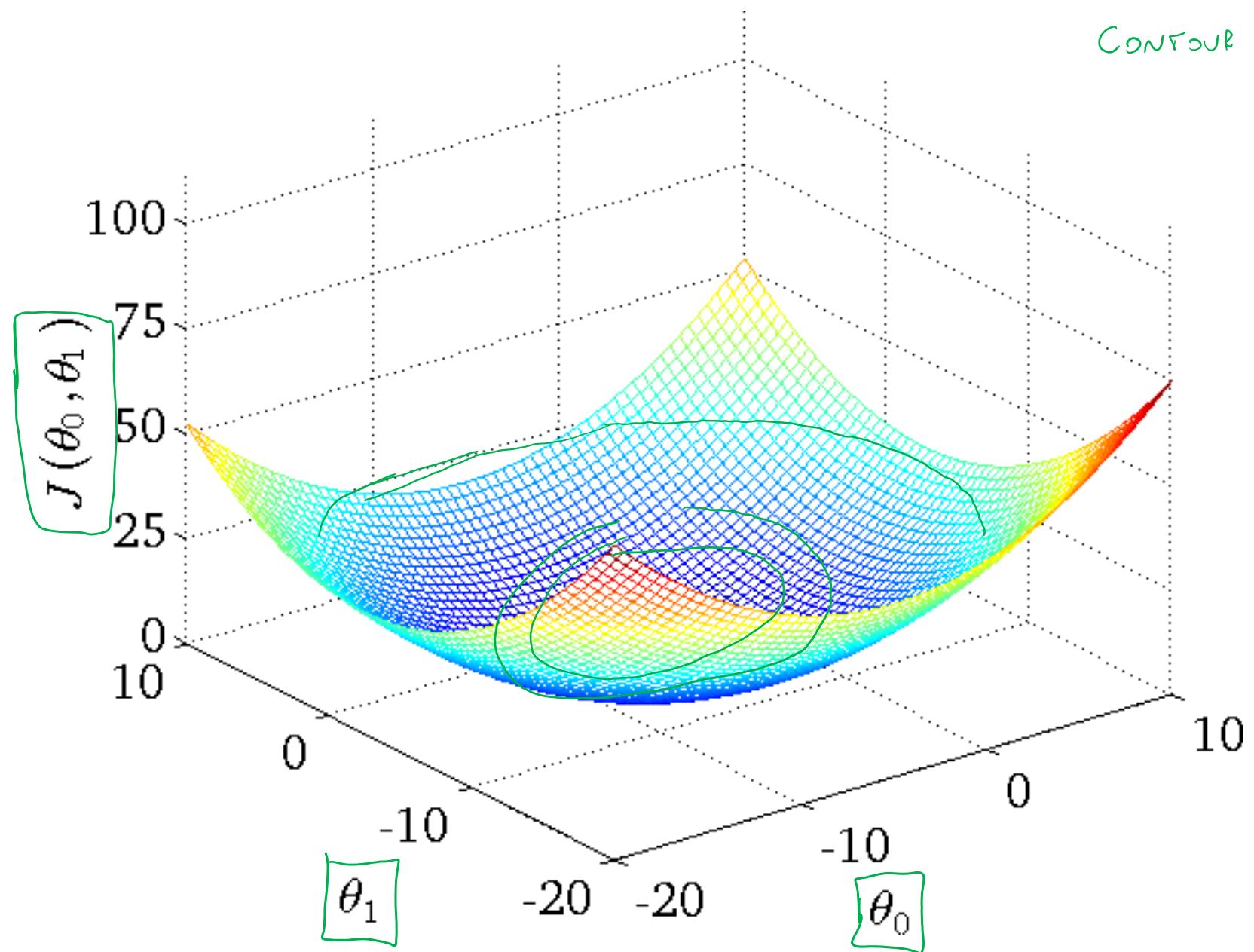
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



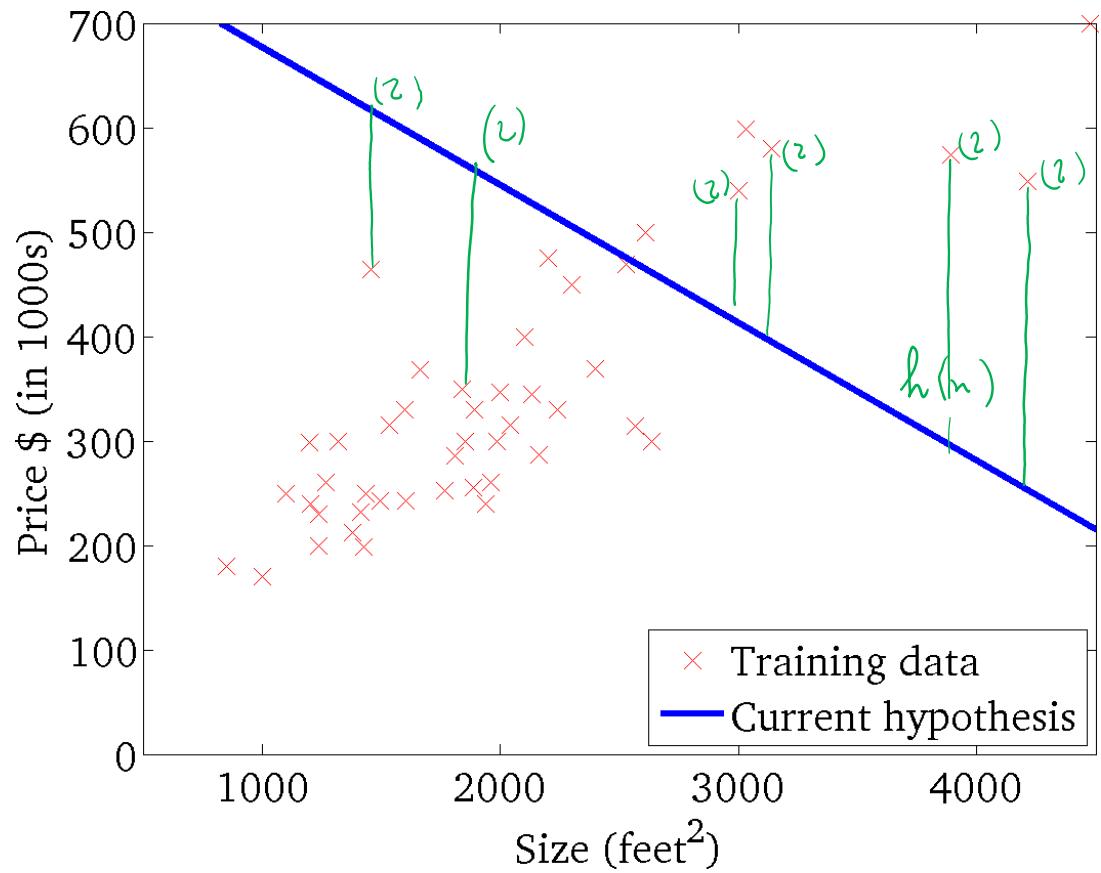
$$(\theta_0, \theta_1)$$

Contour Pors



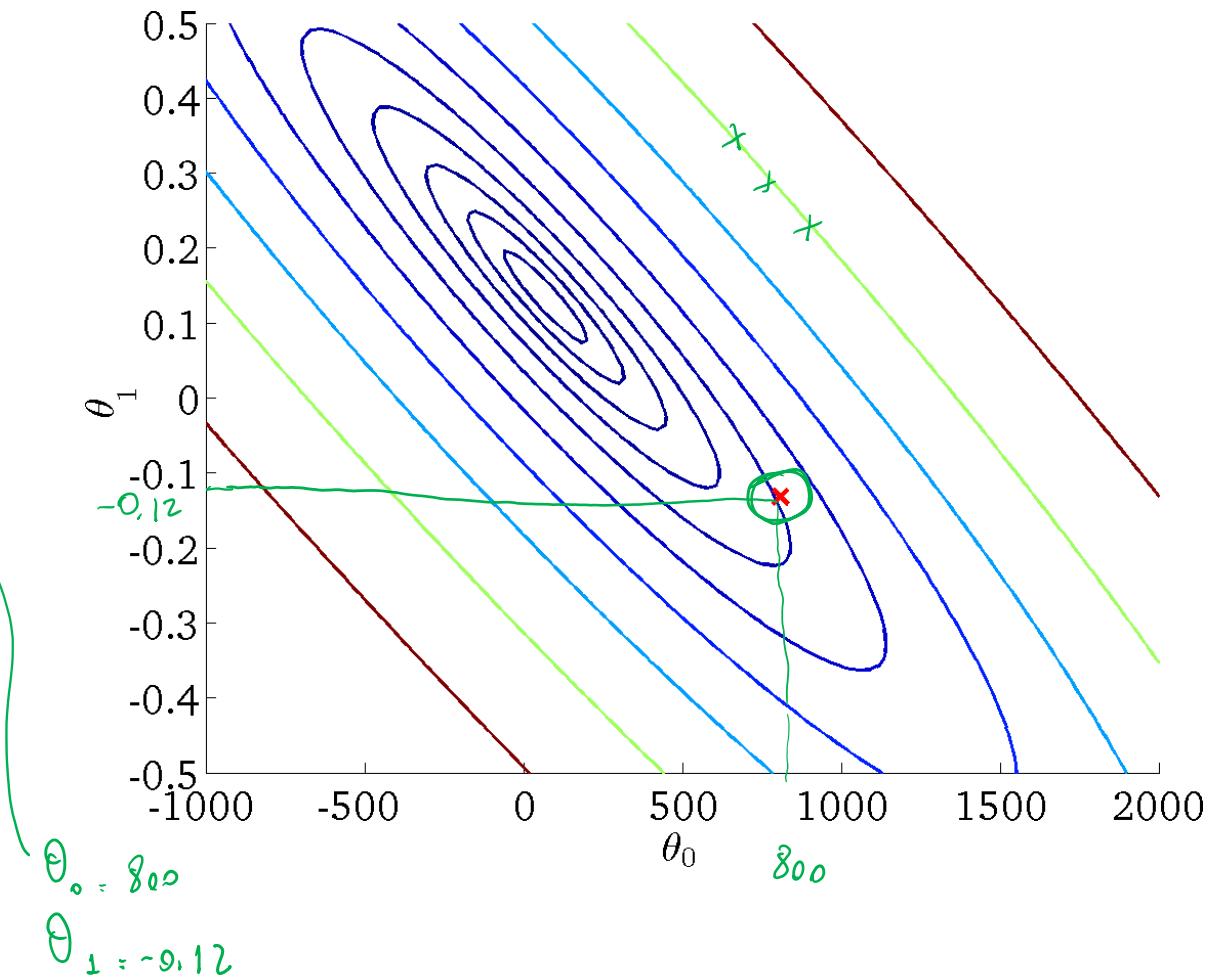
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



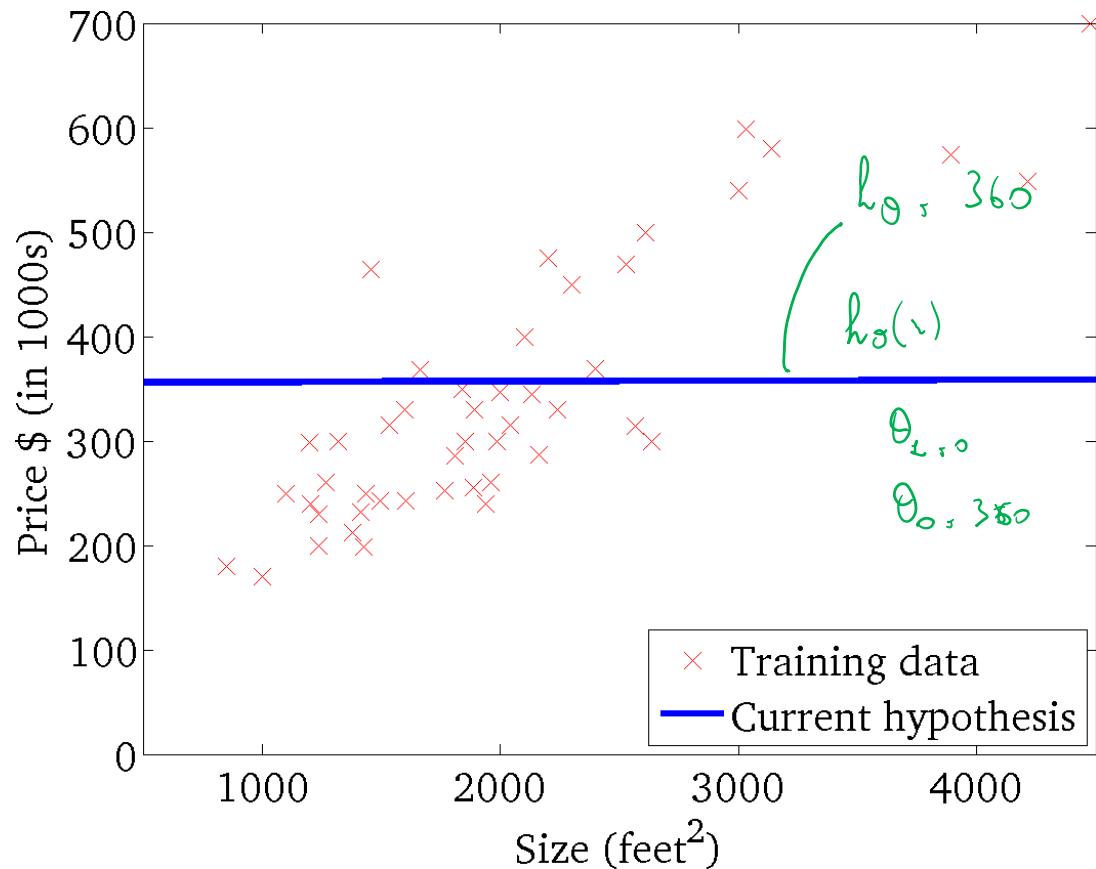
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



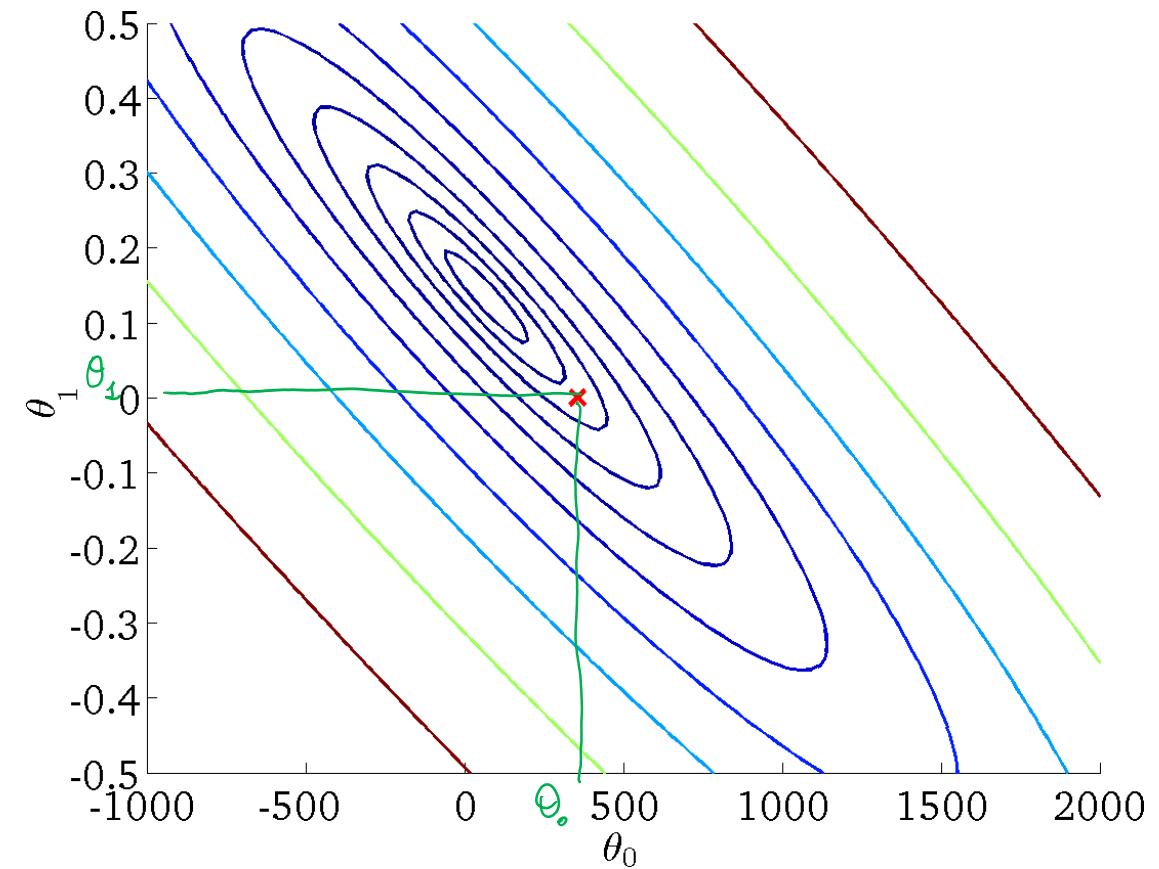
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



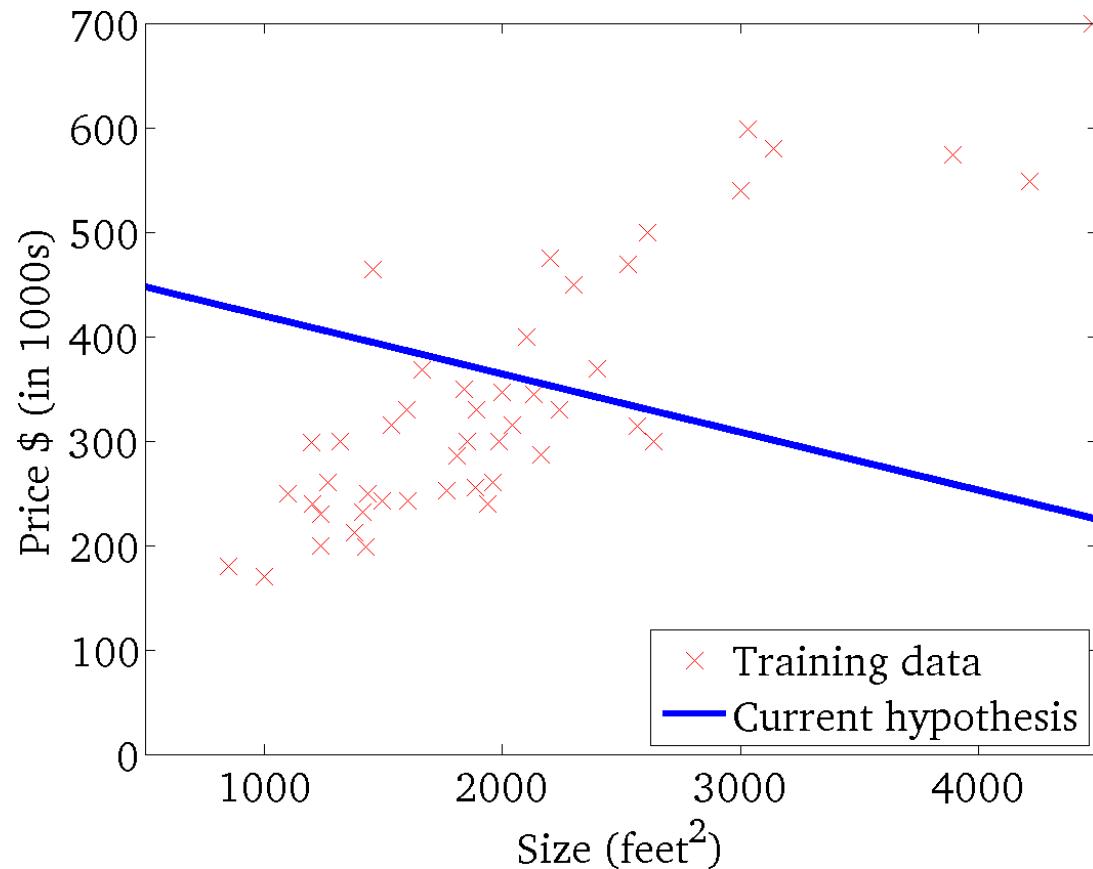
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



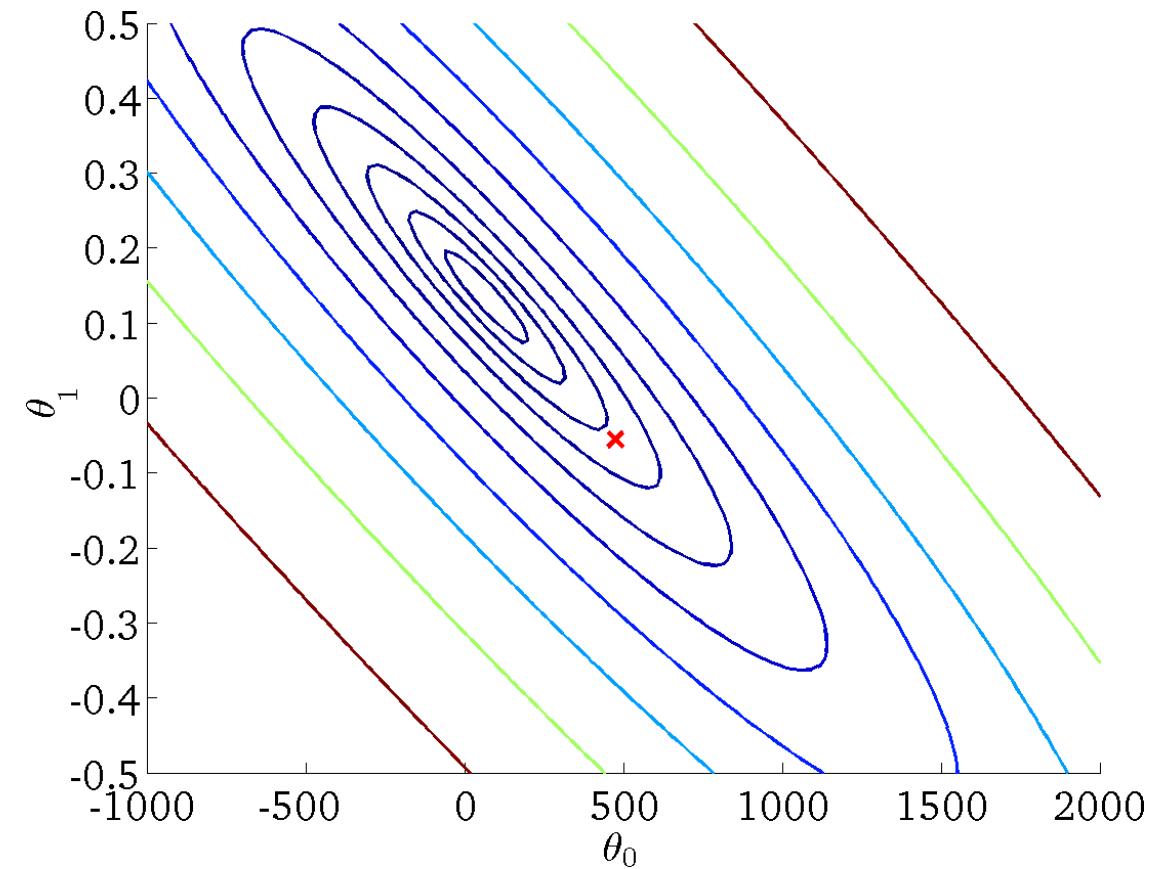
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



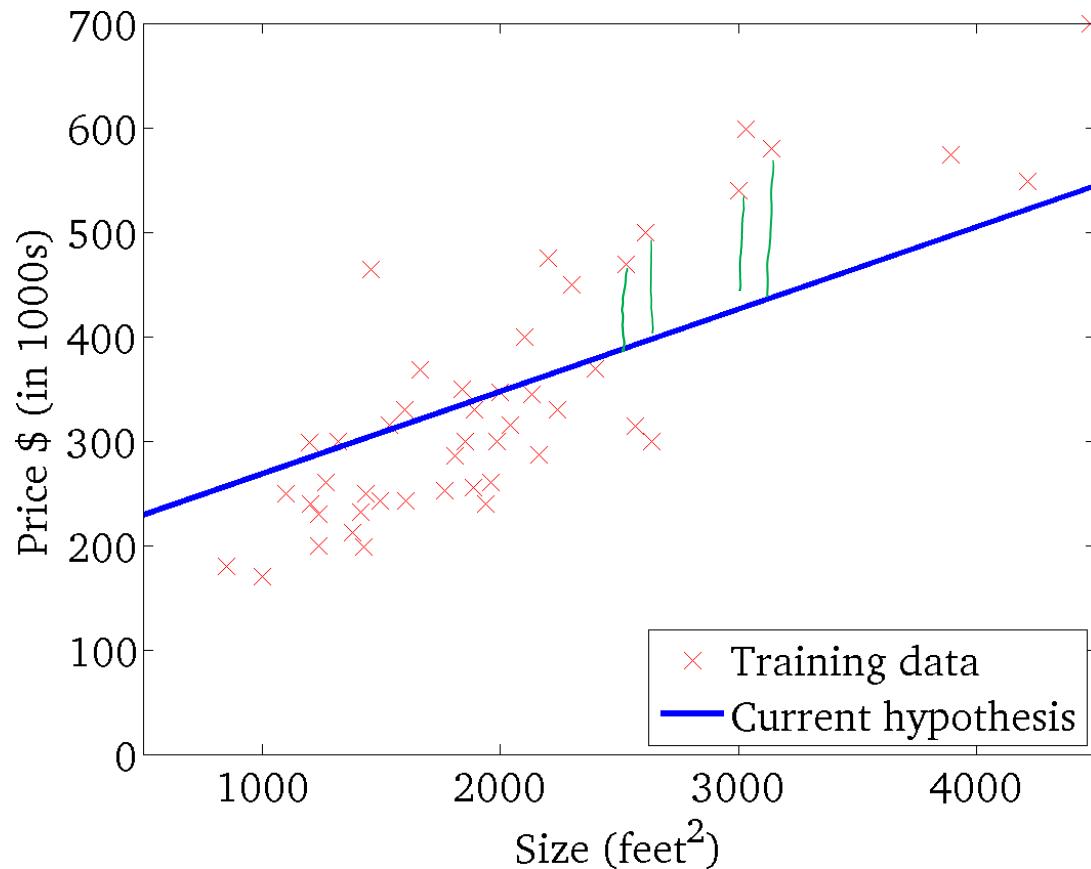
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



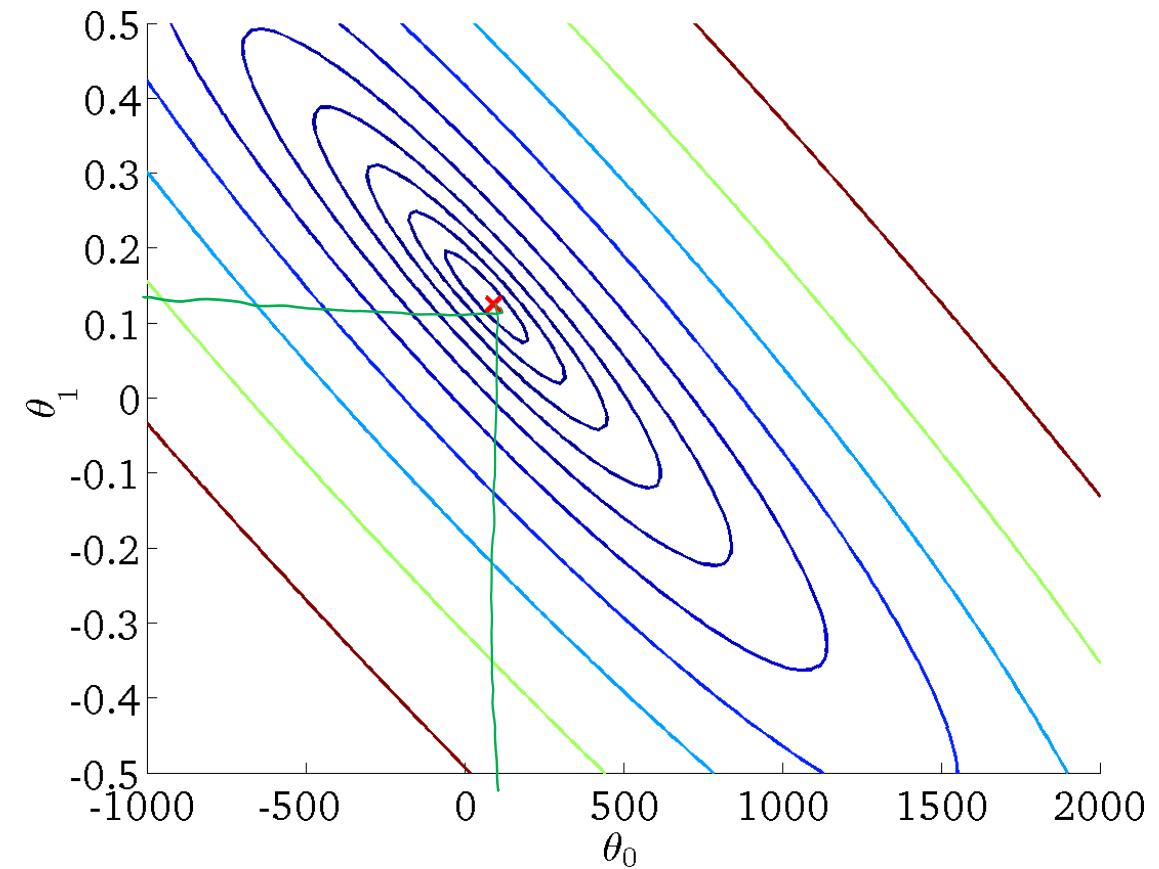
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

(function of the parameters θ_0, θ_1)



Outline

- Linear regression
 - ▶ One or multiple variables
 - ▶ Cost function
 - Gradient descent, incl. batch/stochastic GD
 - Normal equation
 - MSE and Correlation
 - Locally-Weighted Regression
 - Probabilistic Interpretation of Least Squares
-

Gradient Descent

Have some function $J(\theta_0, \theta_1)$ $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

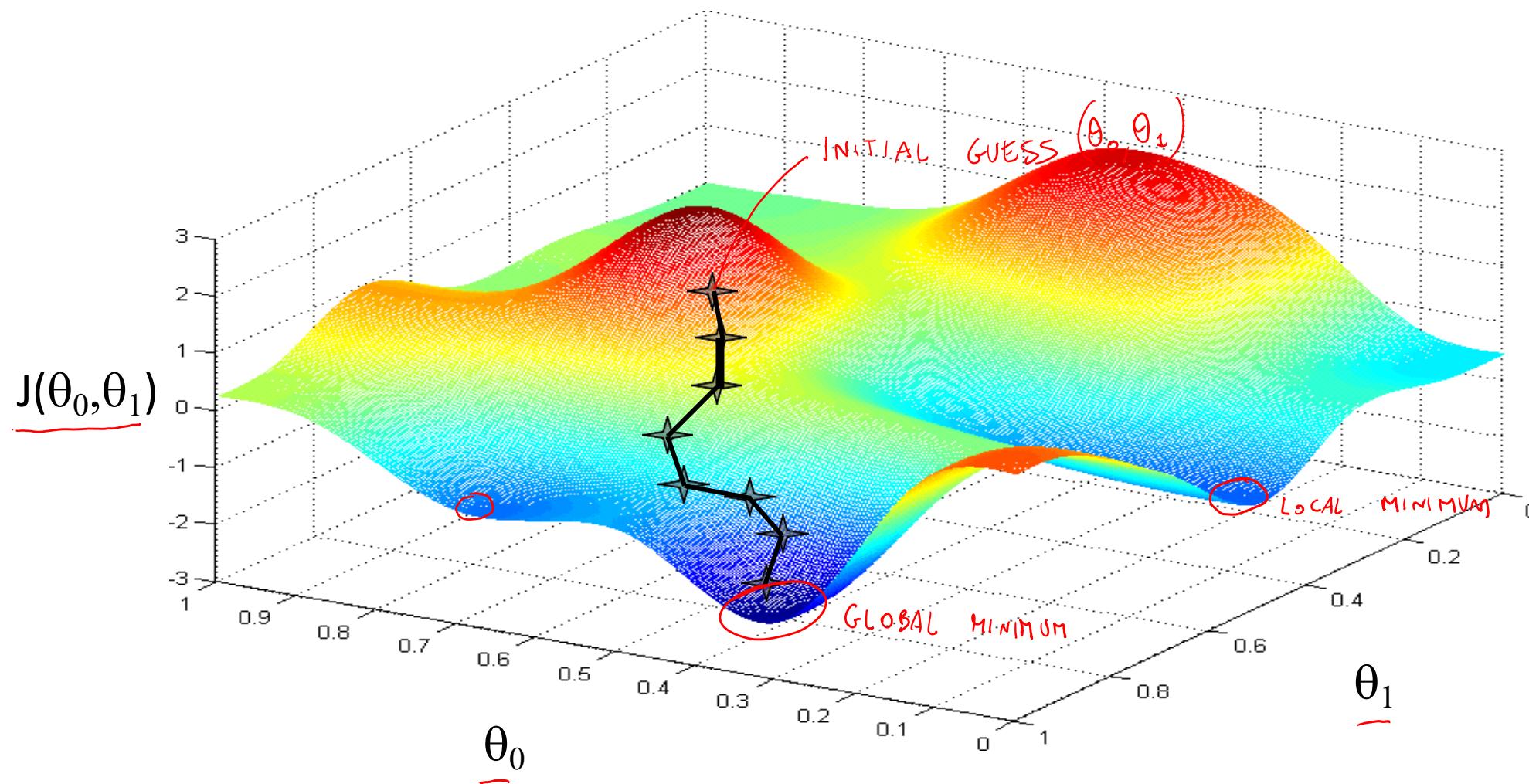
Want $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$ $\min_{\theta_0, \theta_1, \dots, \theta_n} J(\theta_0, \theta_1, \dots, \theta_n)$

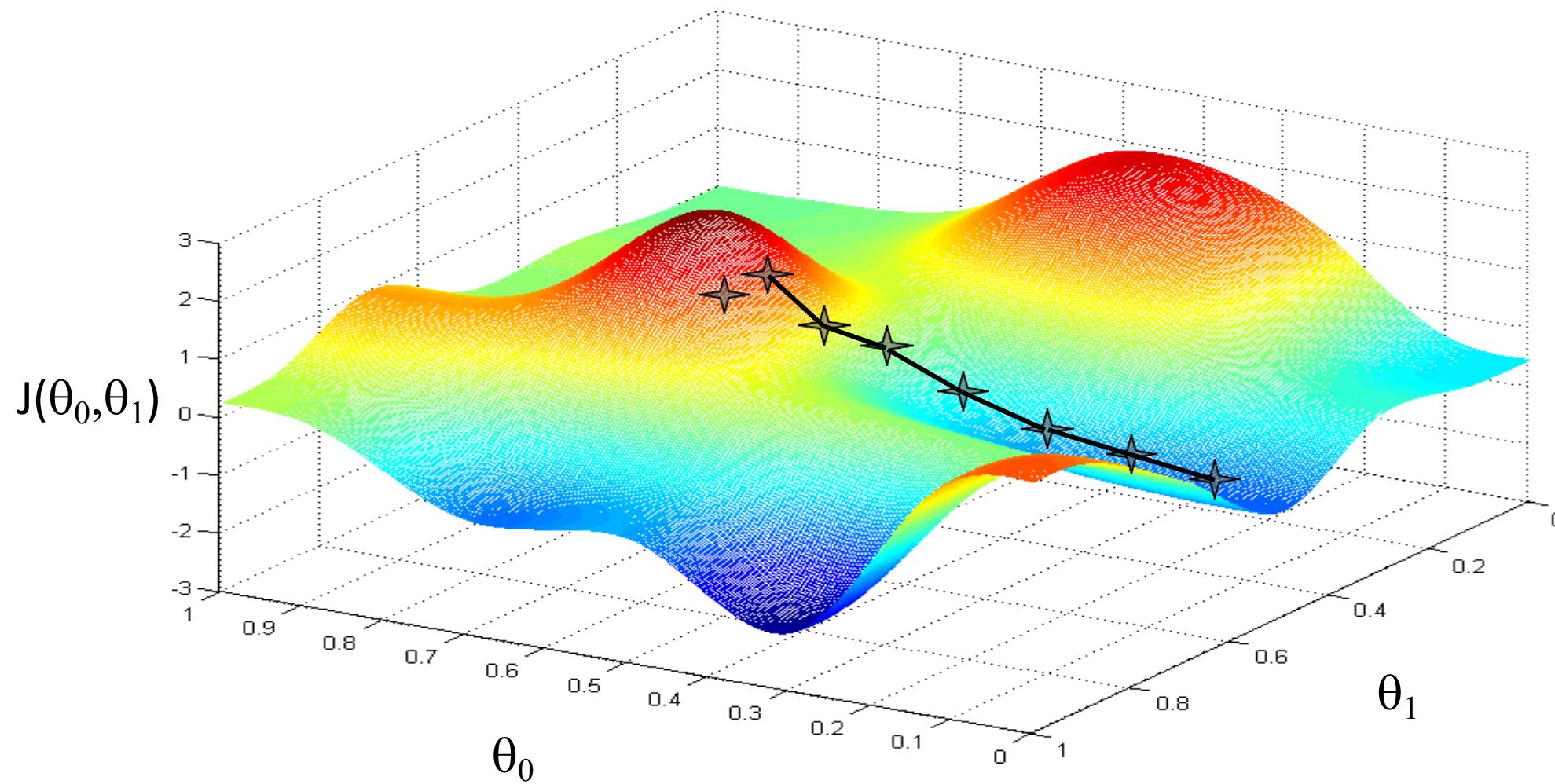
$$y = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots = \theta^T \mathbf{x}$$

$\left(\begin{matrix} x_0 = 1 \end{matrix} \right)$

Outline:

- Start with some θ_0, θ_1 (say $\theta_0=0, \theta_1=0$)
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at a minimum





Gradient descent algorithm

```
repeat until convergence {  
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$  }  
    LEARNING RATE  
    SIMULTANEOUSLY  
     $a := b$  ASSIGNMENT  
     $a = b$  ASSERTION
```

Correct: Simultaneous update

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Incorrect:

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\underline{\theta_0} := \text{temp0}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_1 := \text{temp1}$$

Gradient descent algorithm

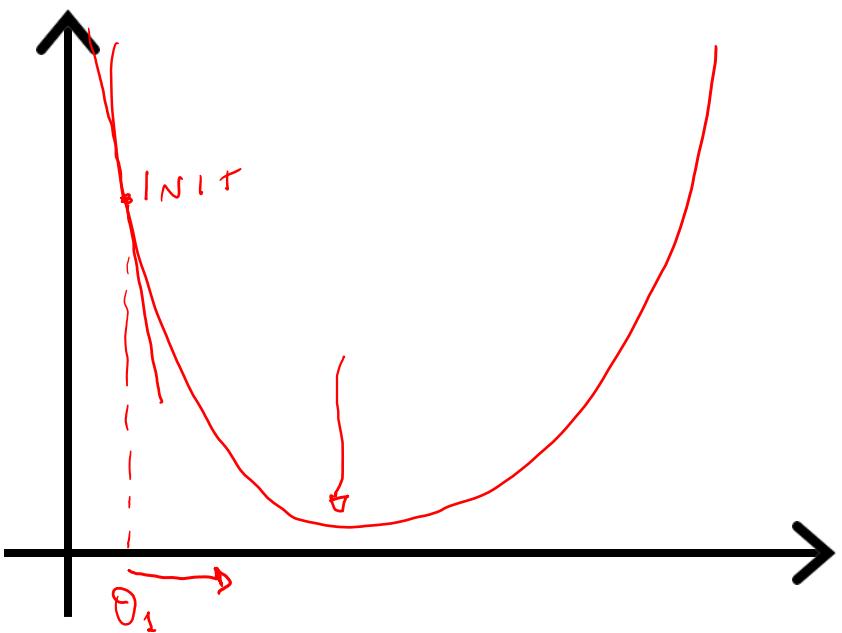
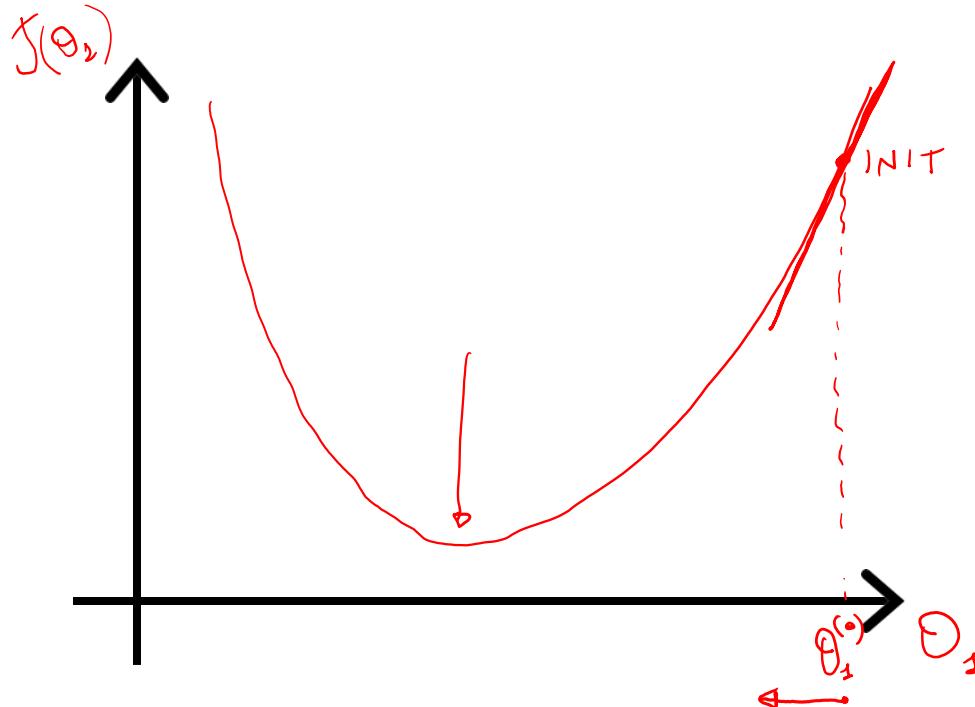
repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \begin{matrix} & \text{(simultaneously update} \\ & j = 0 \text{ and } j = 1) \end{matrix}$$

}

LR DERIVATIVE

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$



$$\Theta_2 \in \mathbb{R}$$

$$\Theta_2^{(t+1)} := \Theta_2^{(t)} - \alpha \underbrace{\frac{\delta}{\delta \Theta_2} J(\Theta_2^{(t)})}_{\geq 0}$$

$$\Theta_2 := \Theta_2 - \alpha \text{ (positive number)}$$

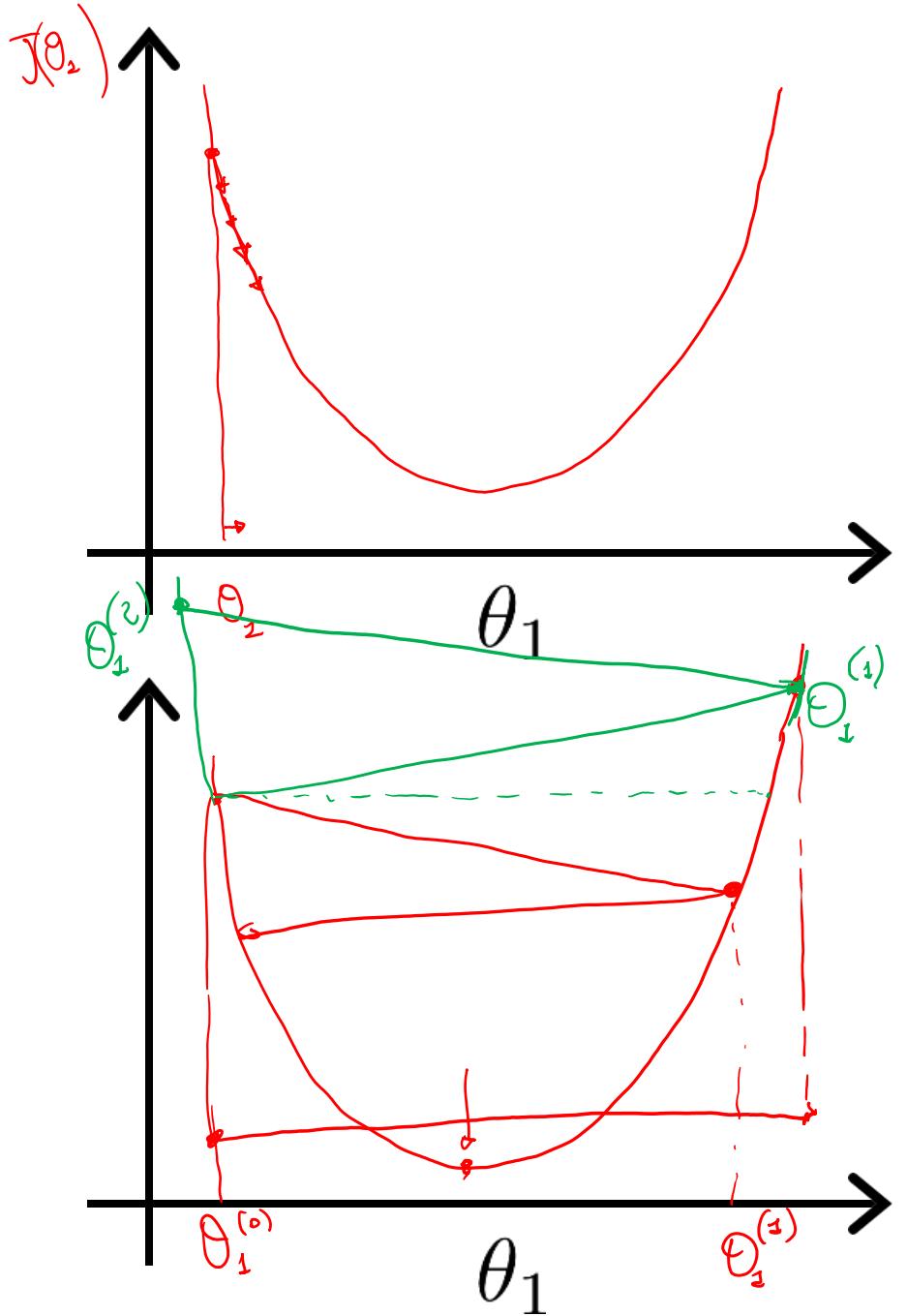
$$\underbrace{\frac{\delta}{\delta \Theta_L} J(\Theta_L)}_{\leq 0}$$

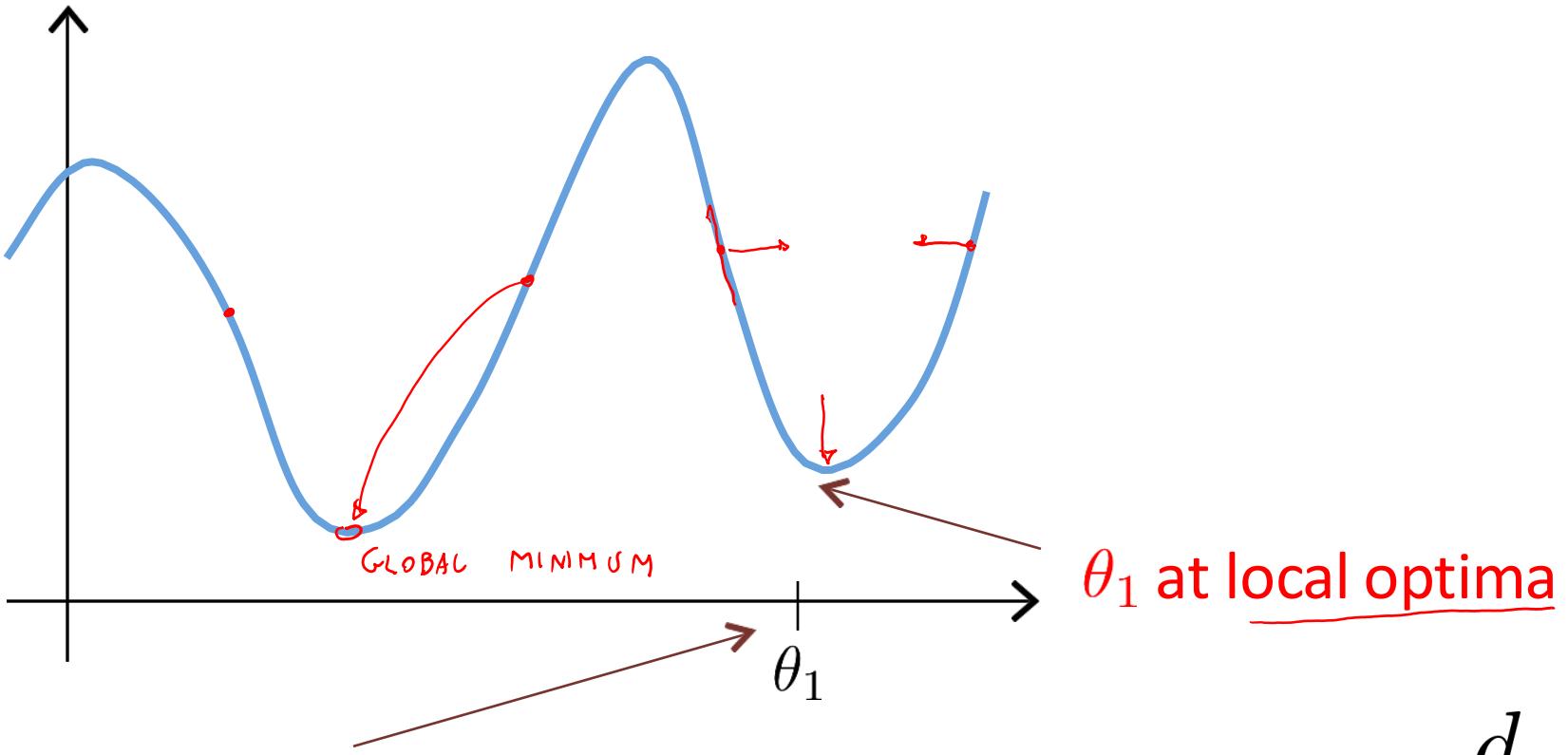
$$\Theta_1 := \Theta_1 - \alpha \text{ (negative number)}$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

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.



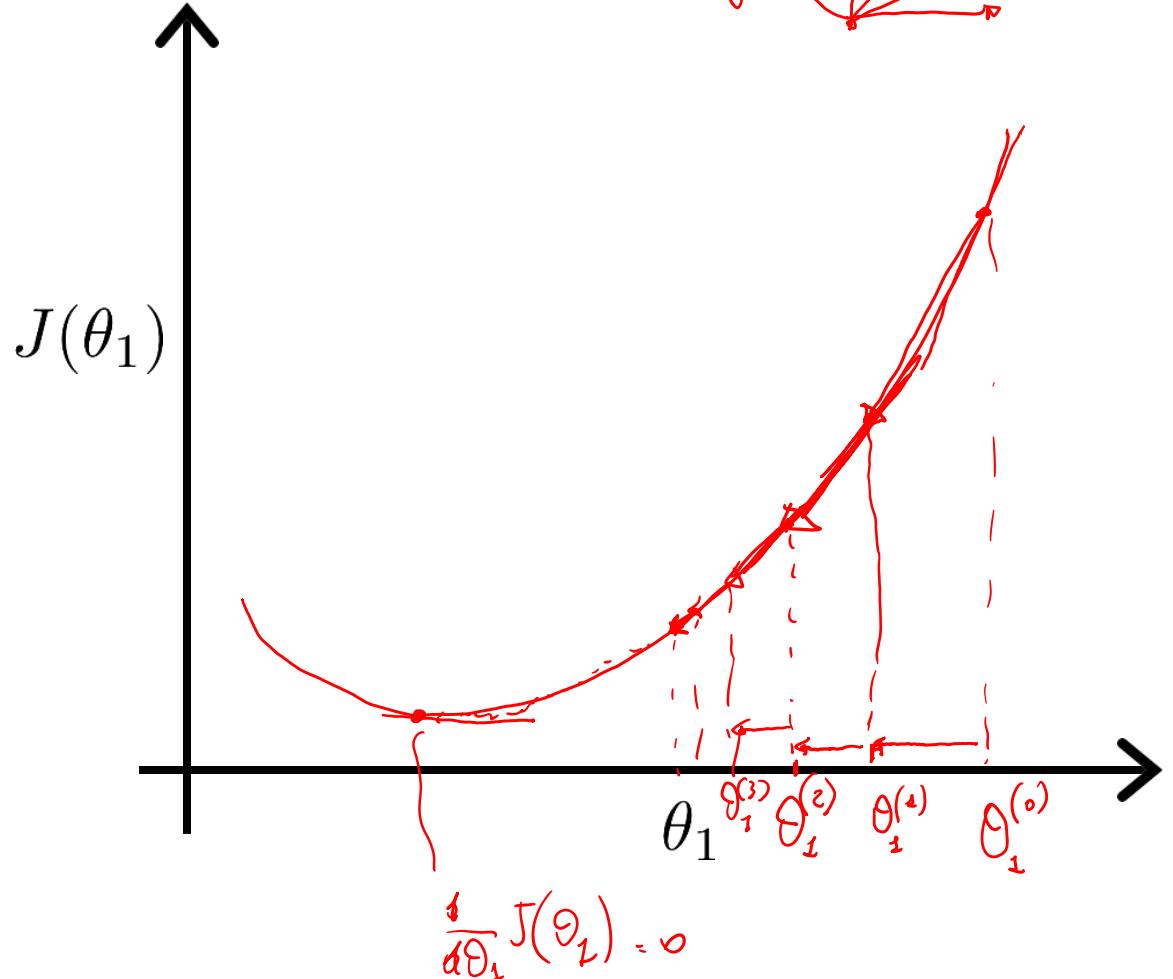


$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

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

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

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



Gradient Descent: for Linear Regression

Gradient descent algorithm

```
repeat until convergence {  
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$   
    (for  $j = 1$  and  $j = 0$ )  
}
```

Linear Regression Model

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

↑
PARAMETRIZED BY

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

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m \left(h_\theta(x^{(i)}) - y^{(i)} \right)^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m \left(\theta_0 + \theta_1 x^{(i)} - y^{(i)} \right)^2\end{aligned}$$

↑

$$j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(h_\theta(x^{(i)}) - y^{(i)} \right) \cdot 1$$

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

Gradient descent algorithm

repeat until convergence {

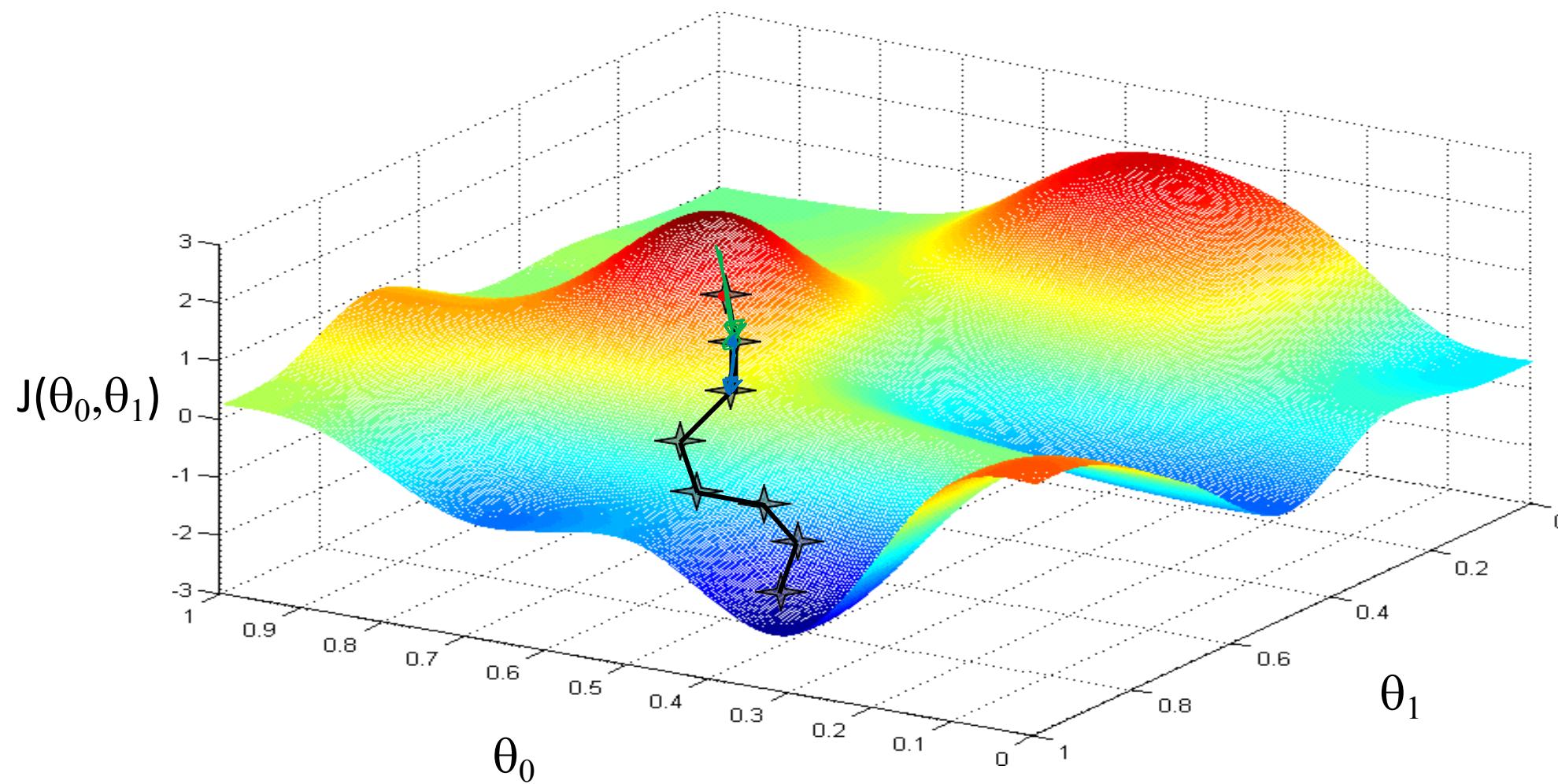
$$\theta_0 := \theta_0 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right]$$
$$\theta_1 := \theta_1 - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \right]$$

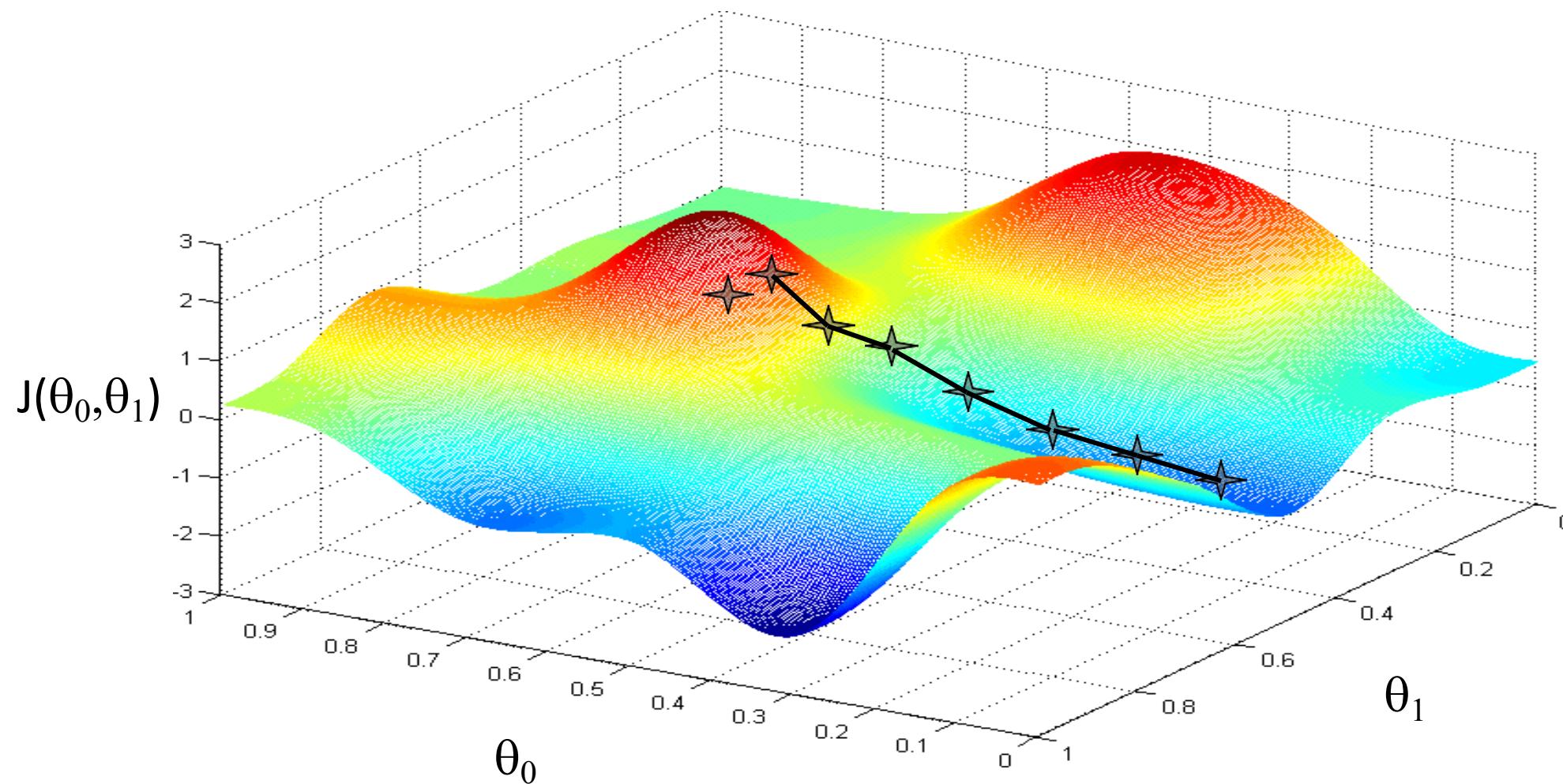
}

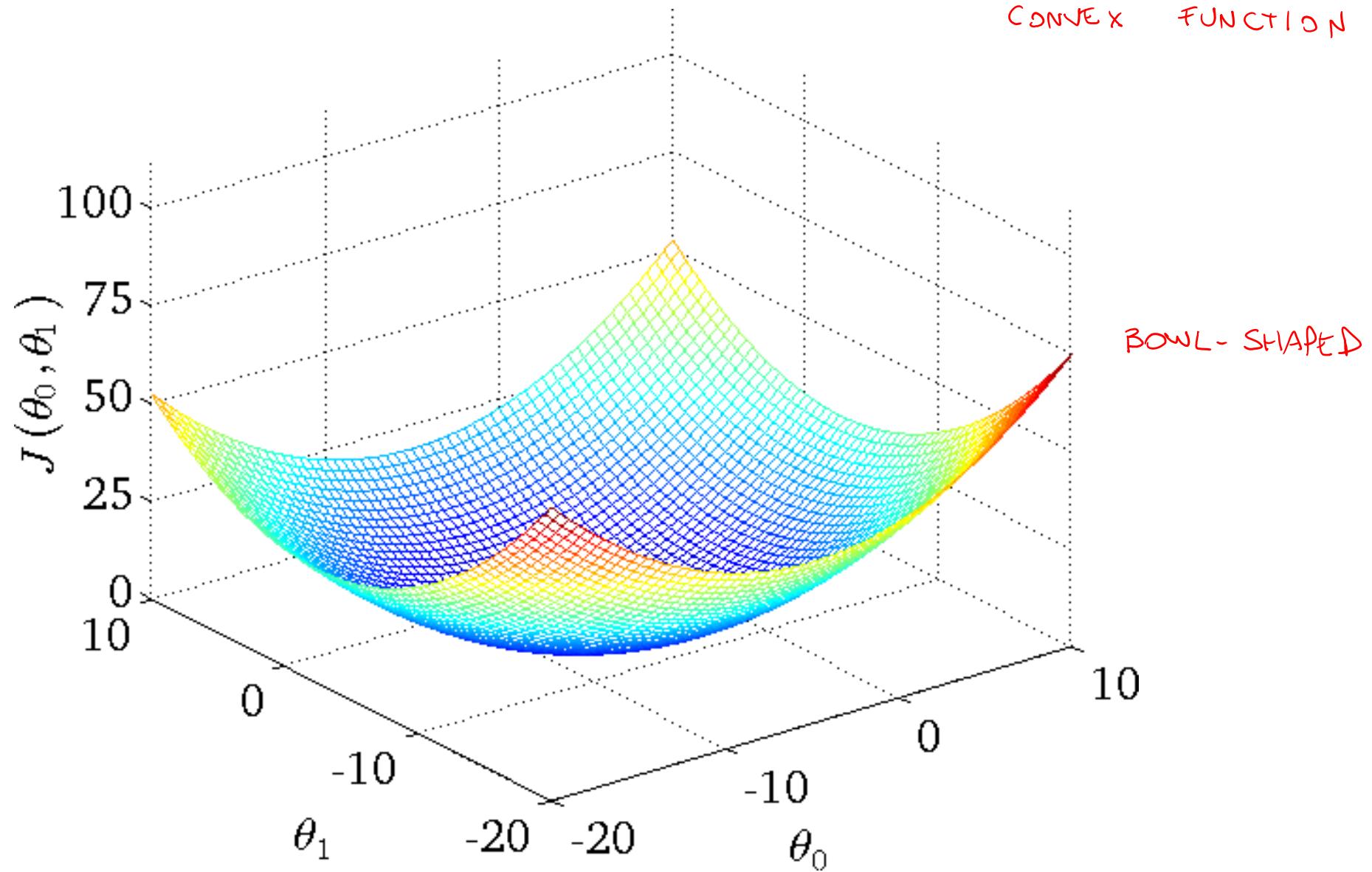
$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

update
 θ_0 and θ_1
simultaneously

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

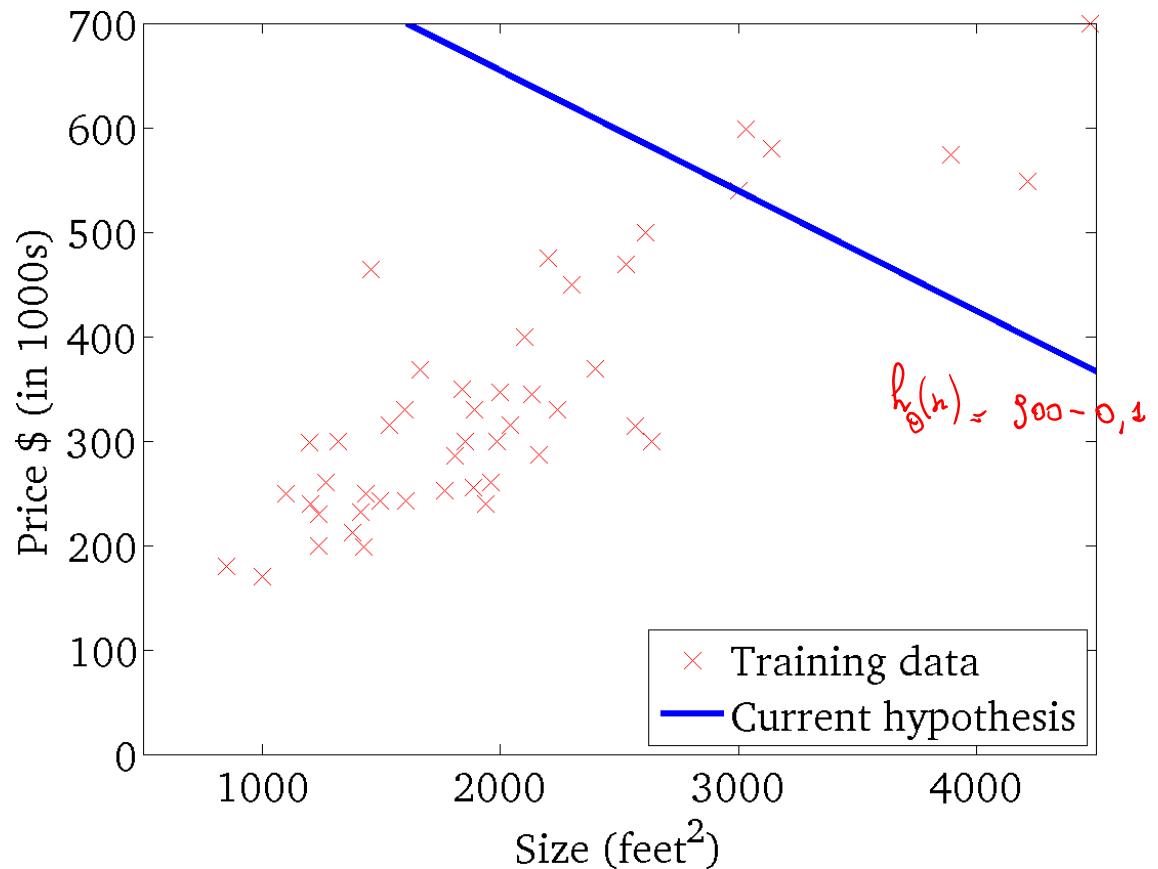






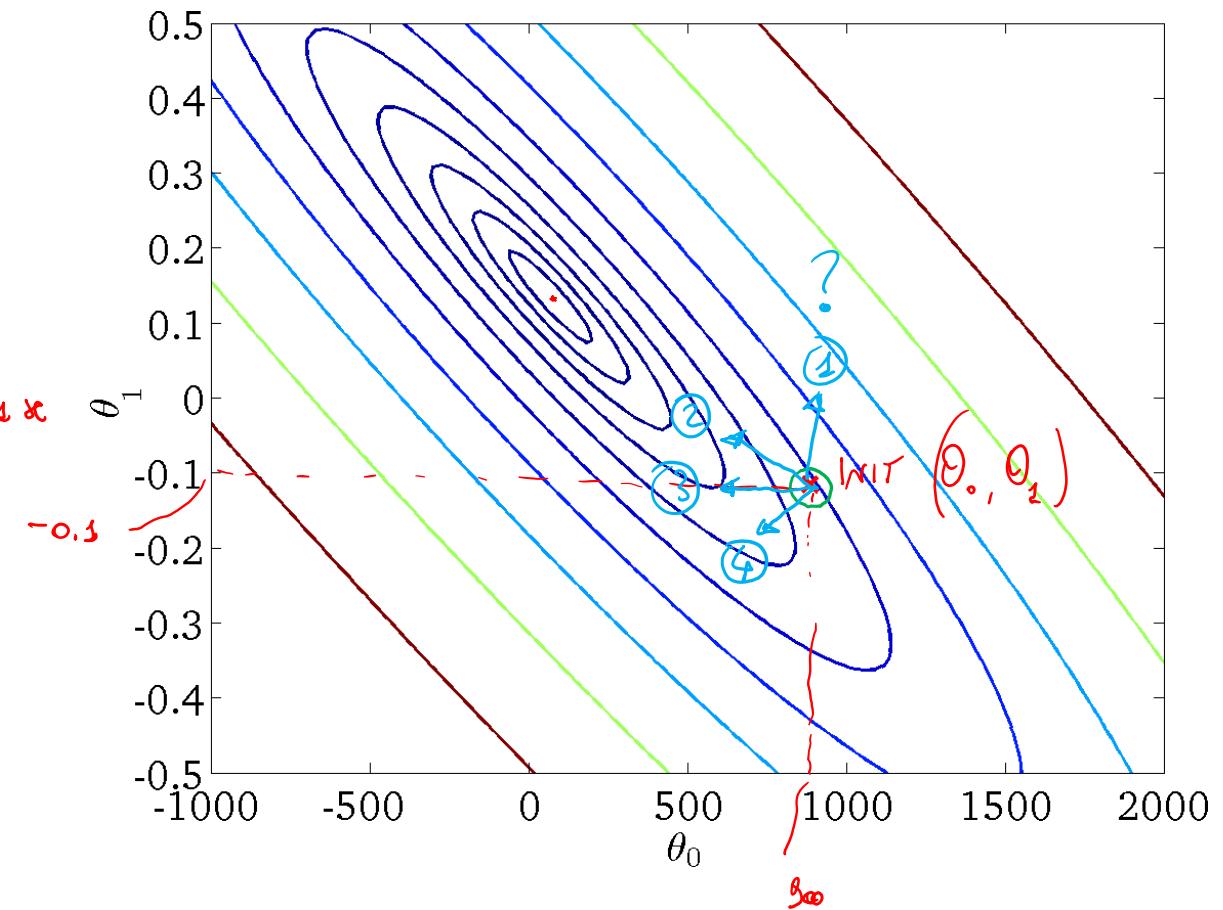
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



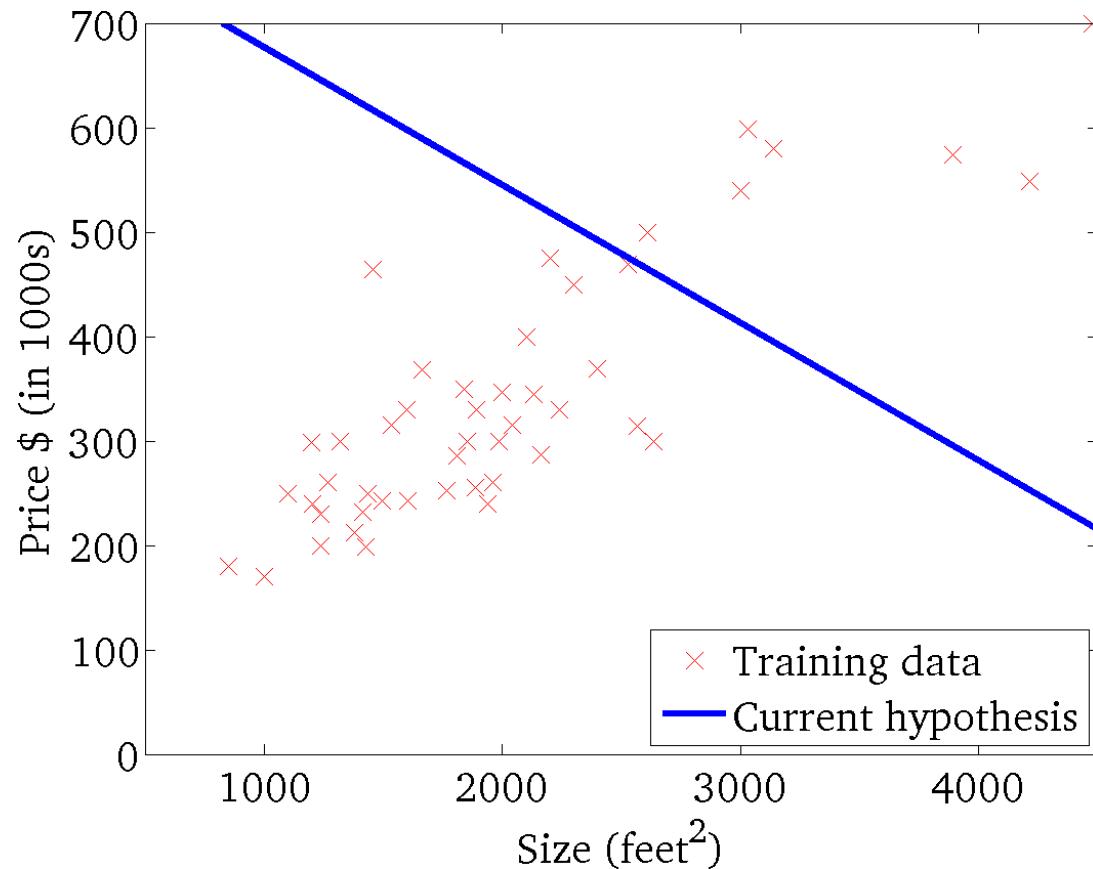
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



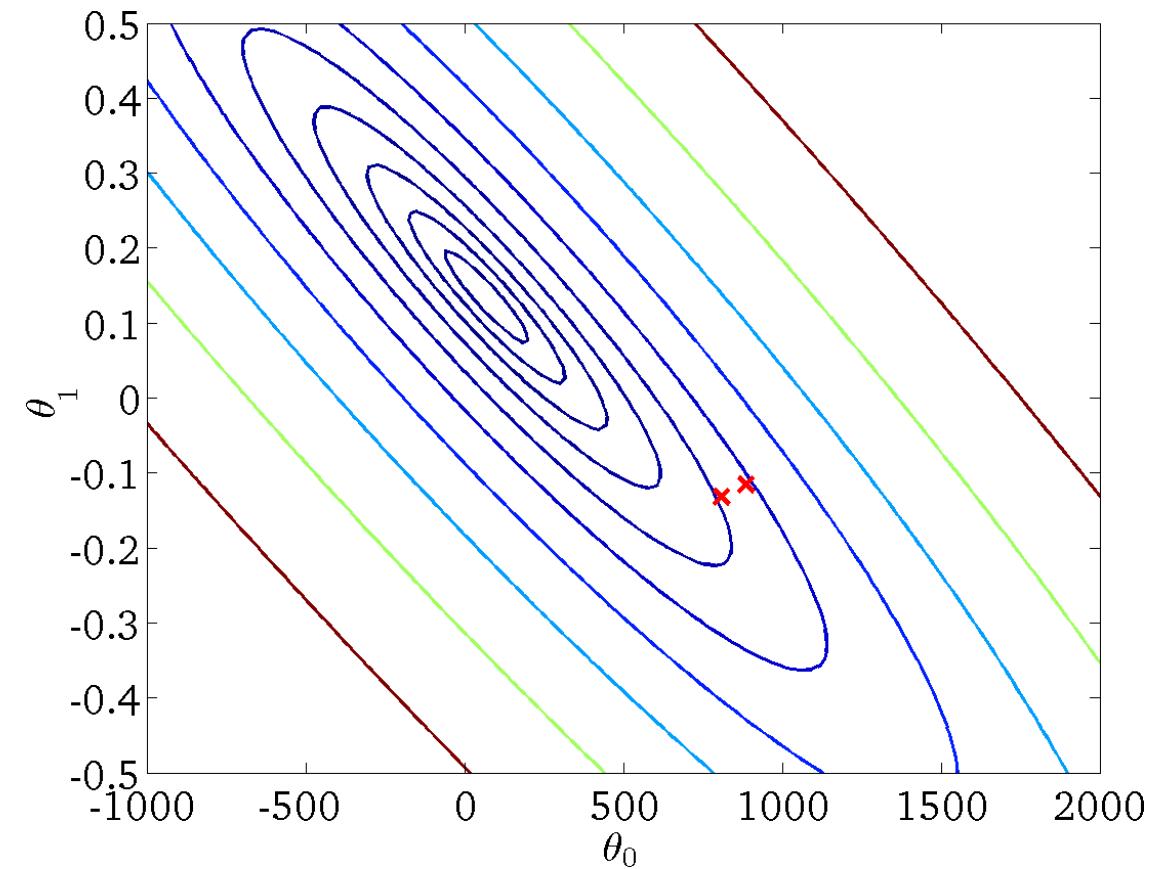
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



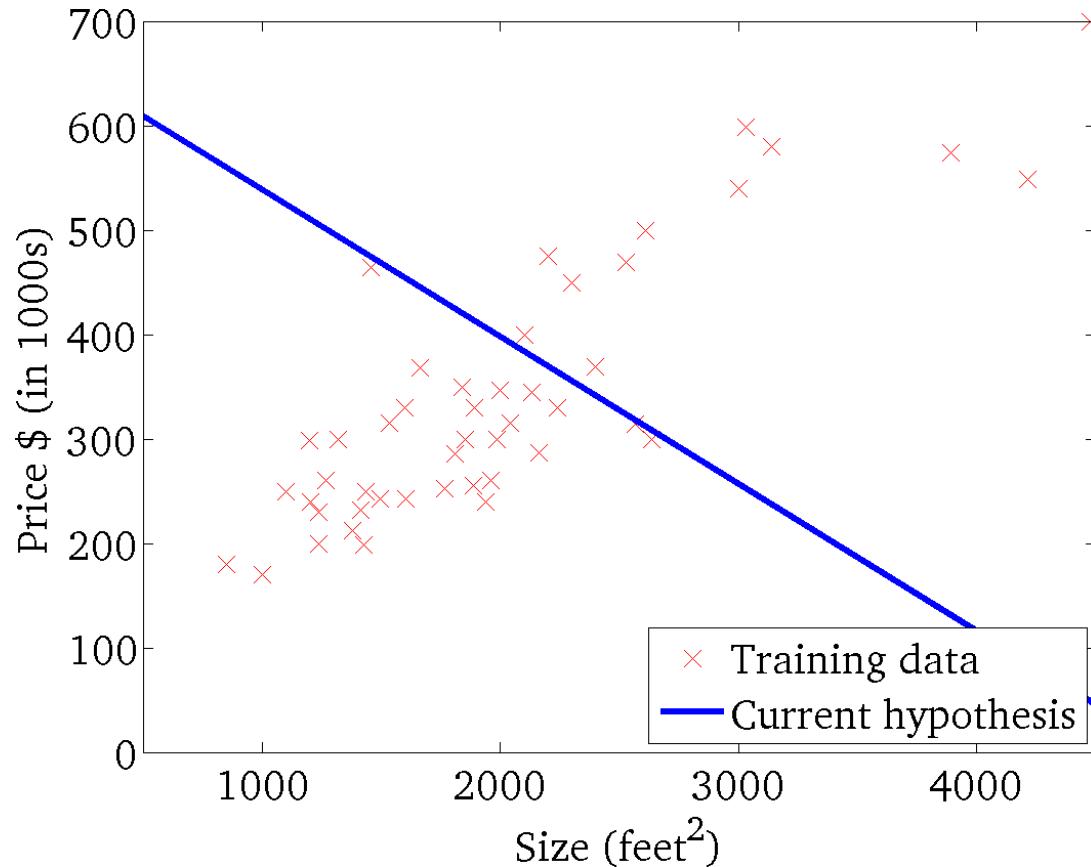
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



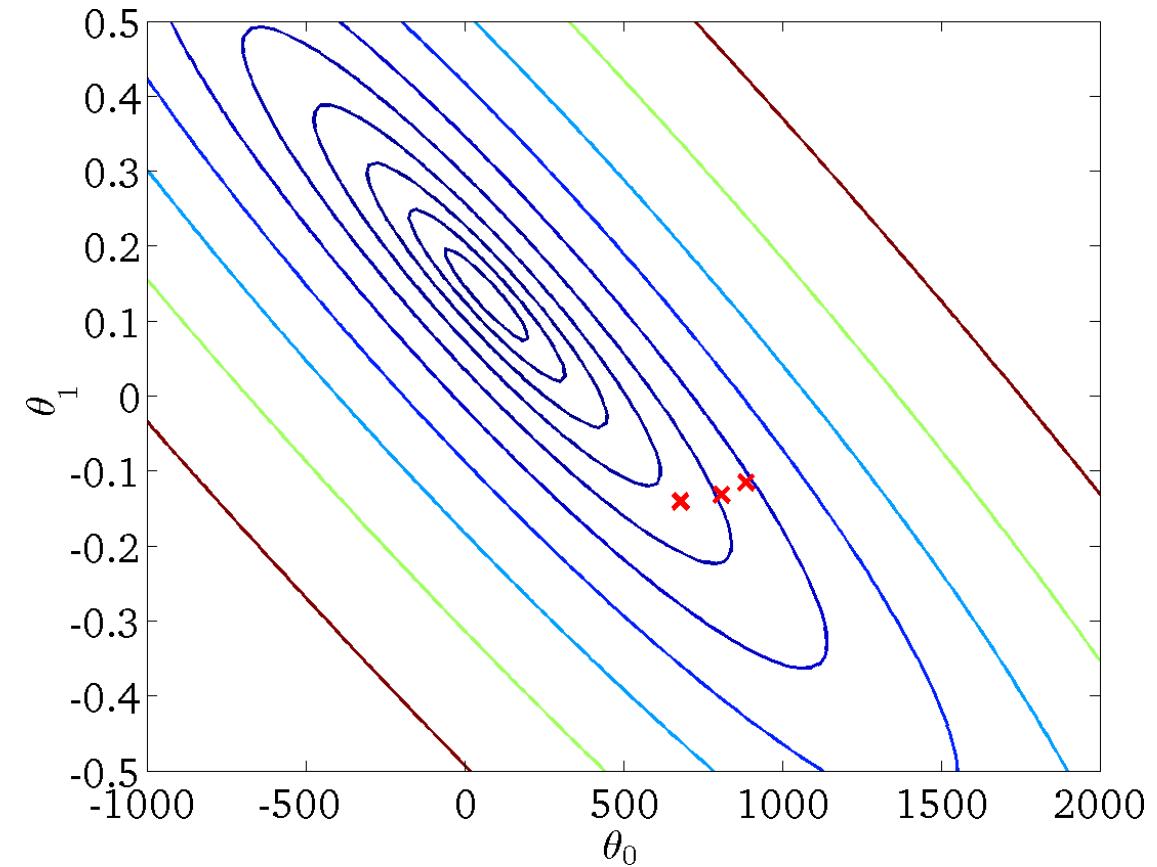
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



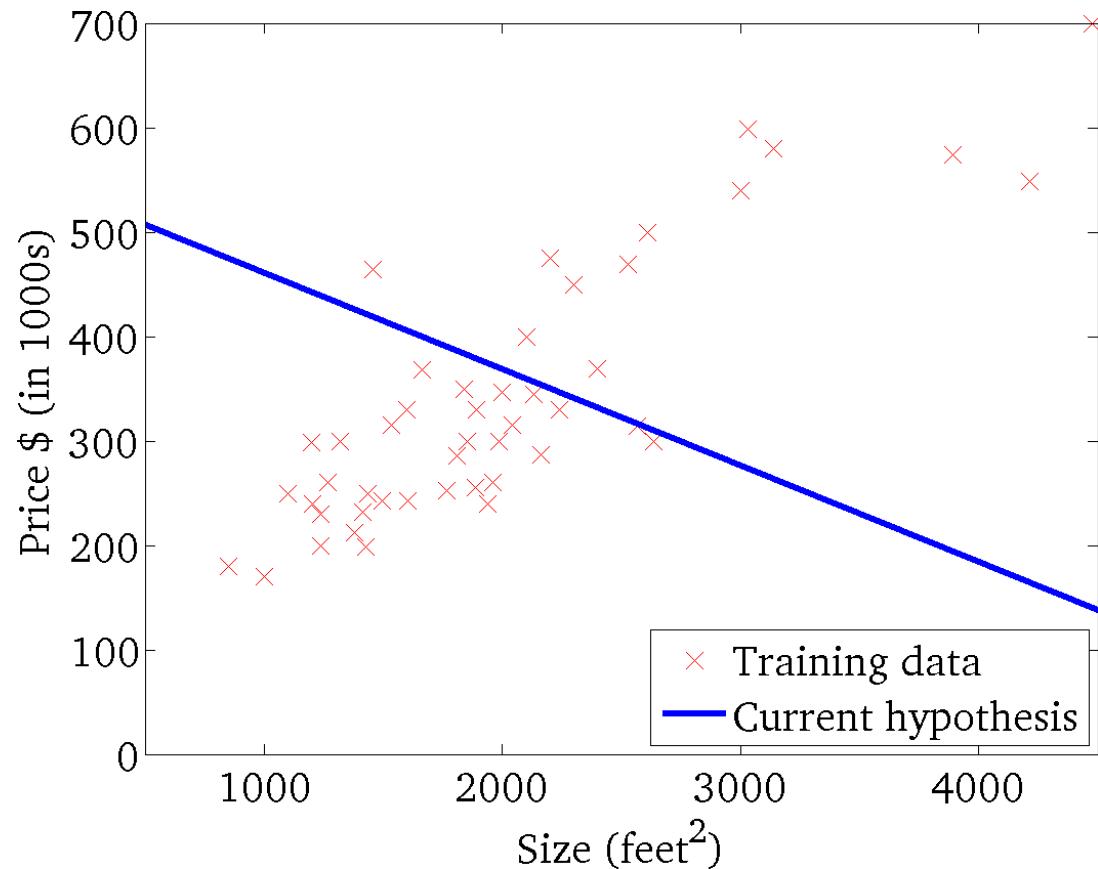
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



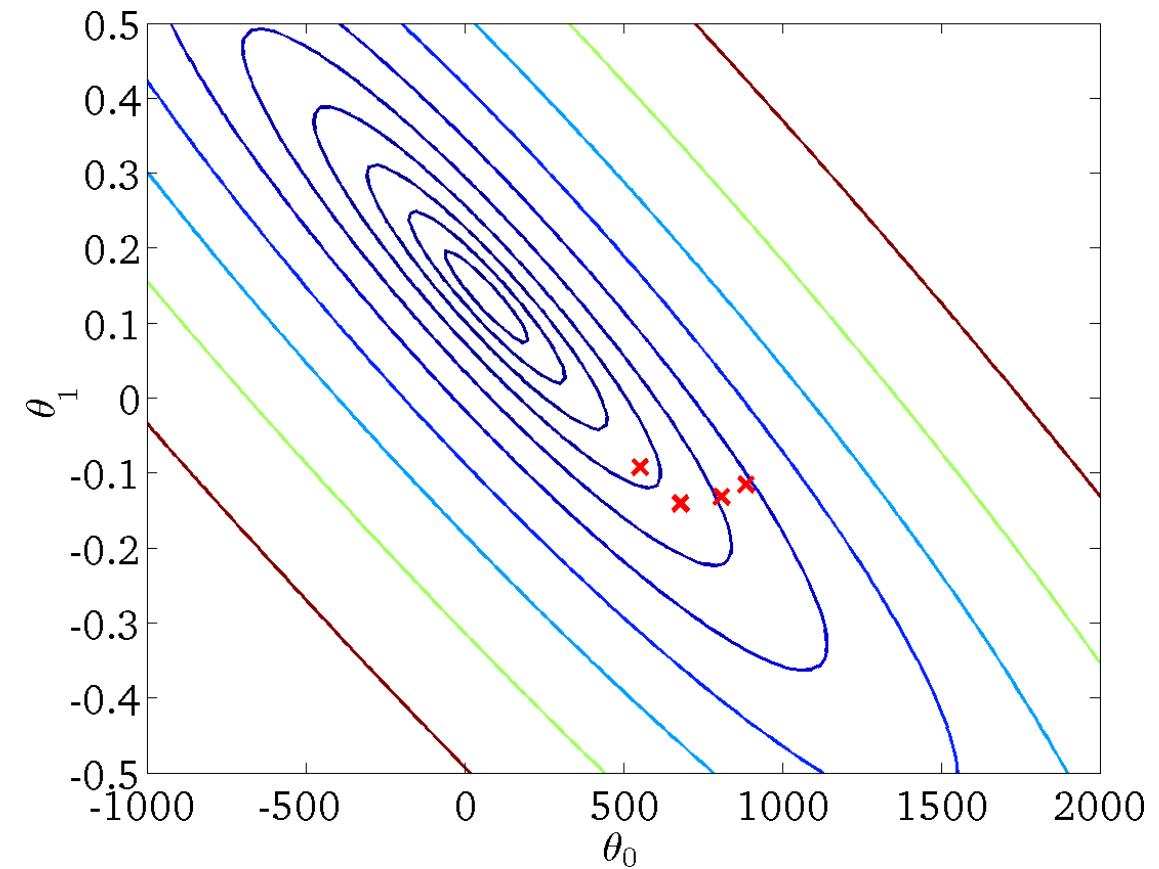
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



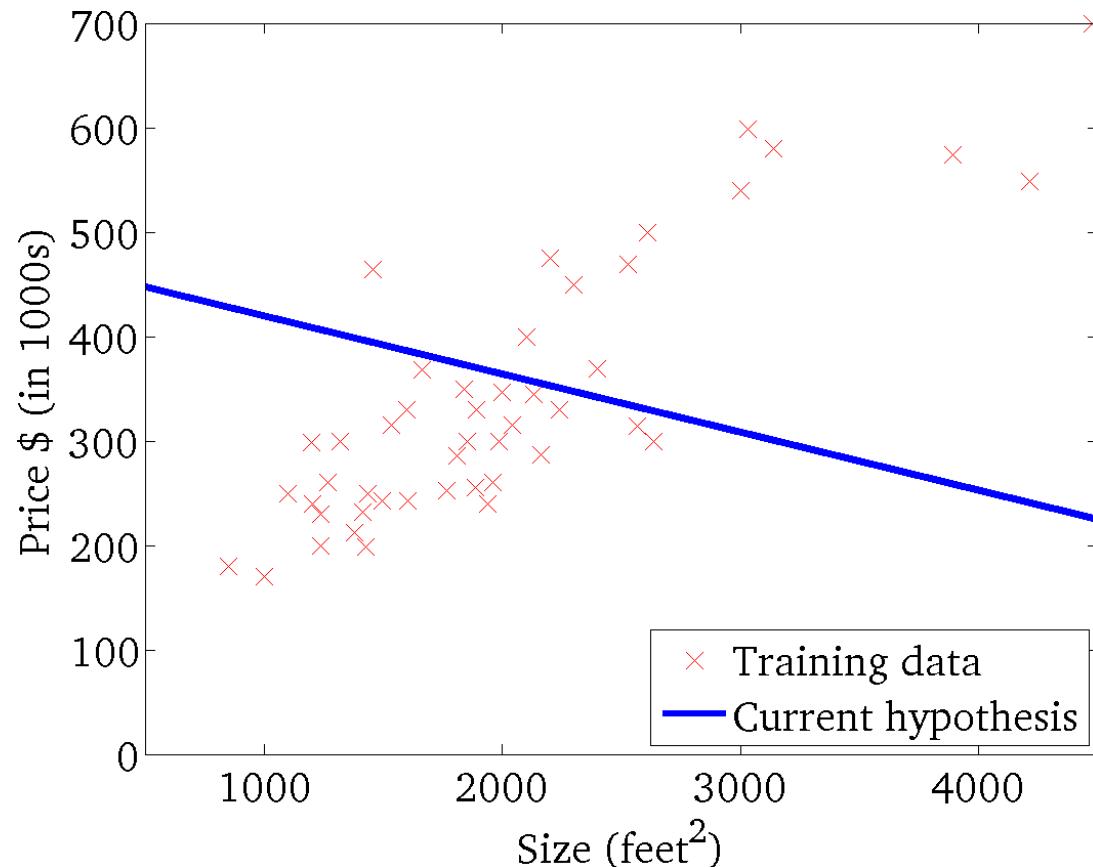
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



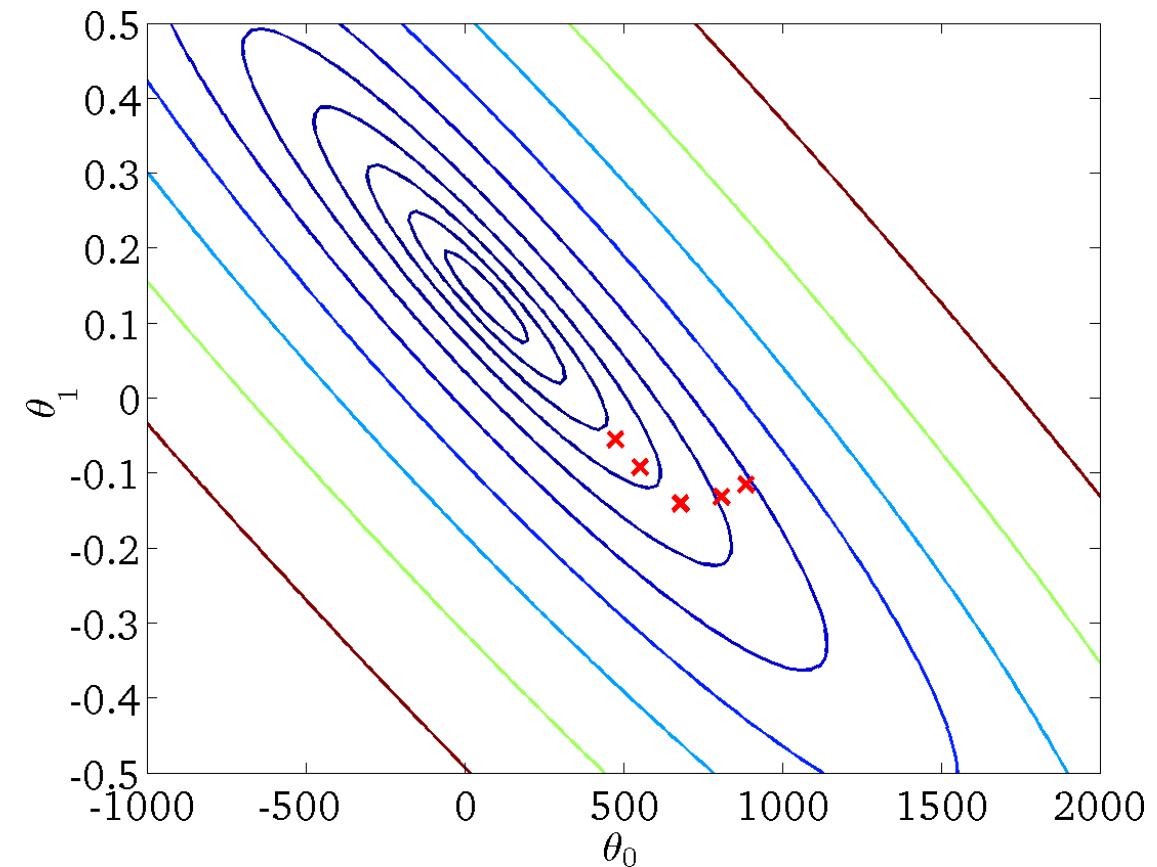
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



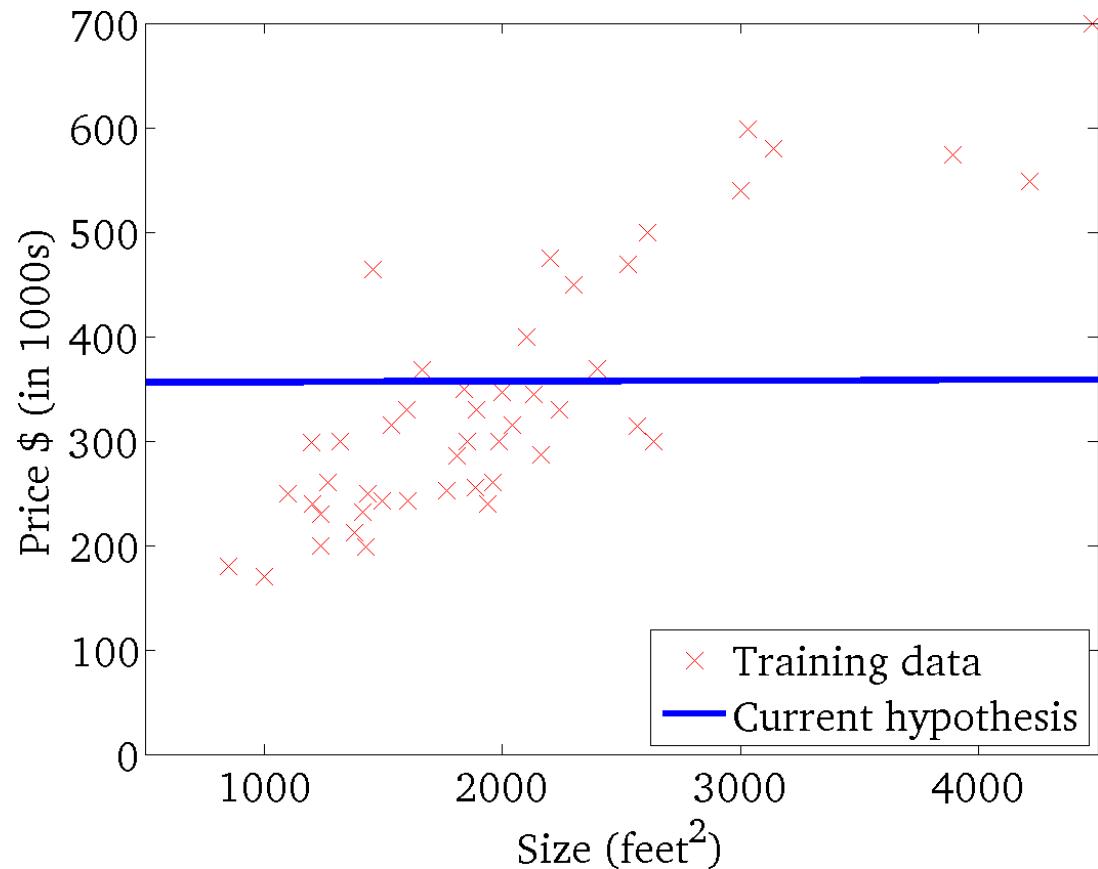
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



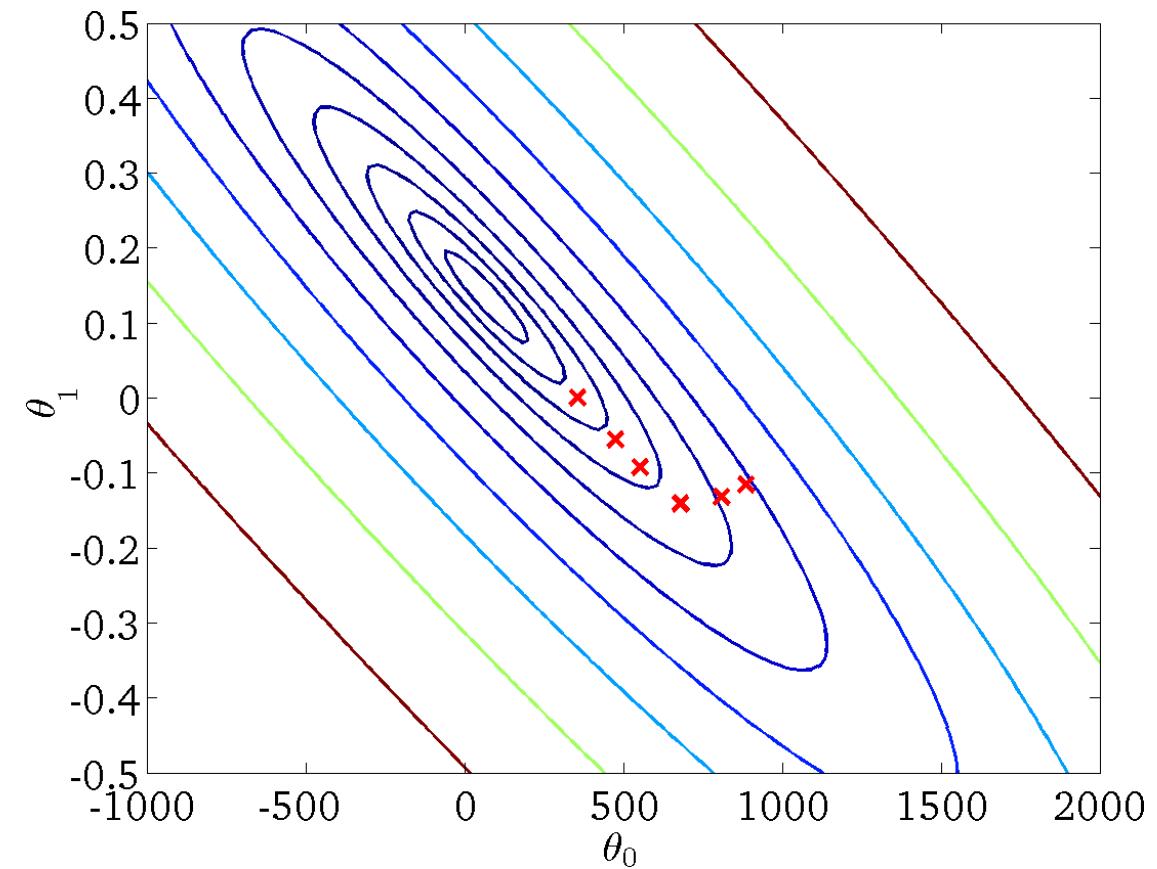
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



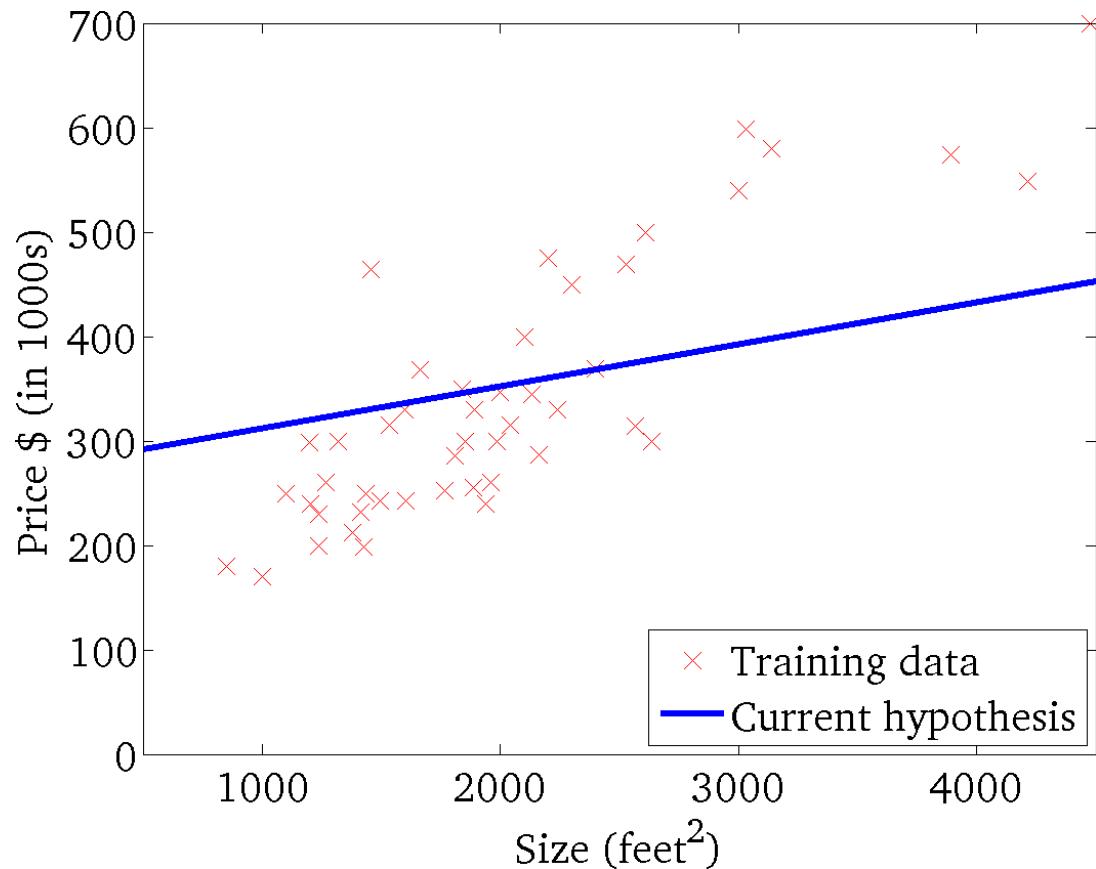
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



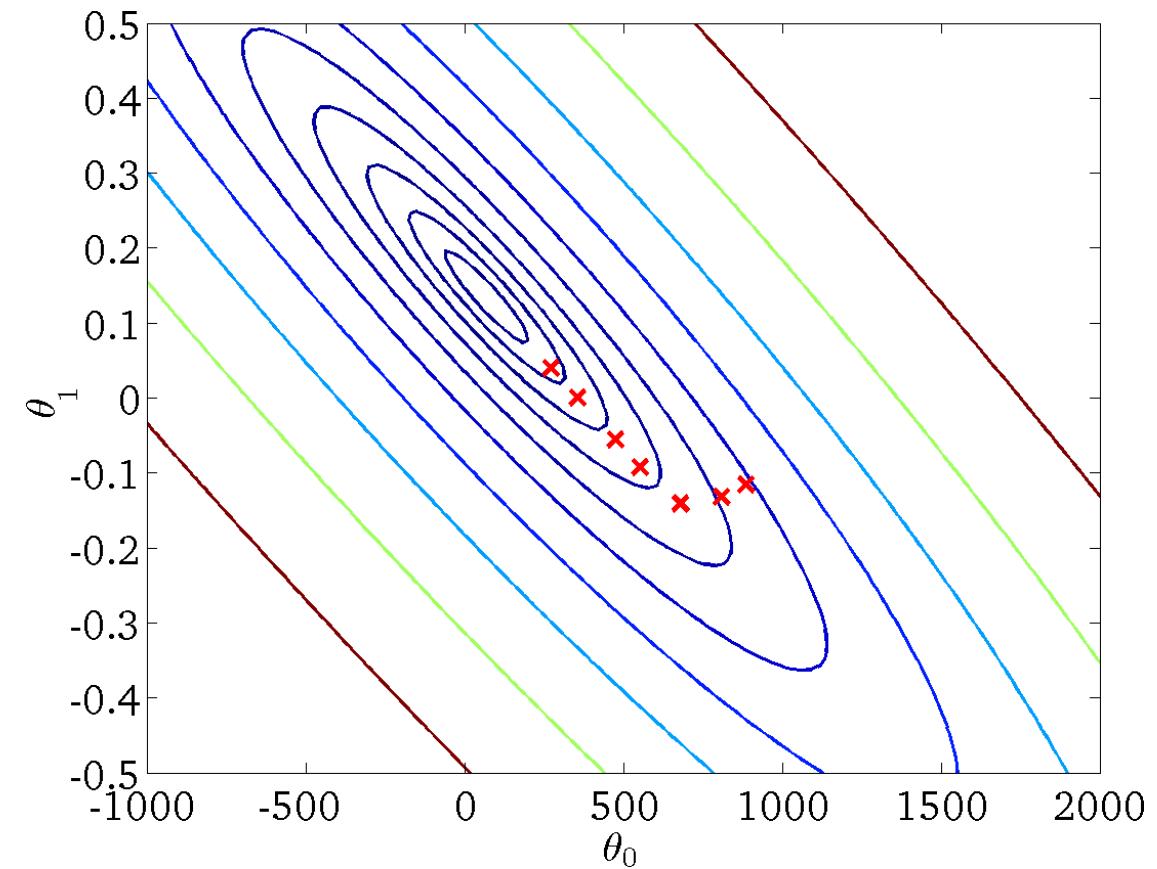
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



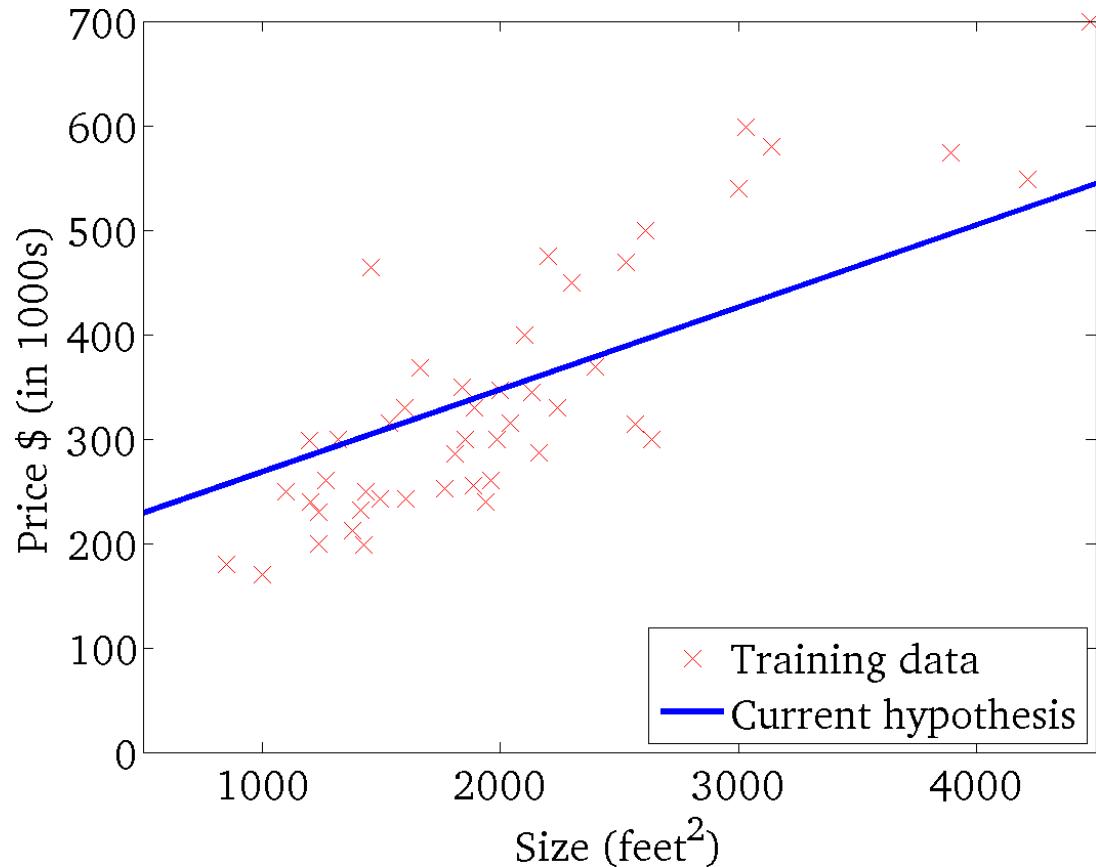
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



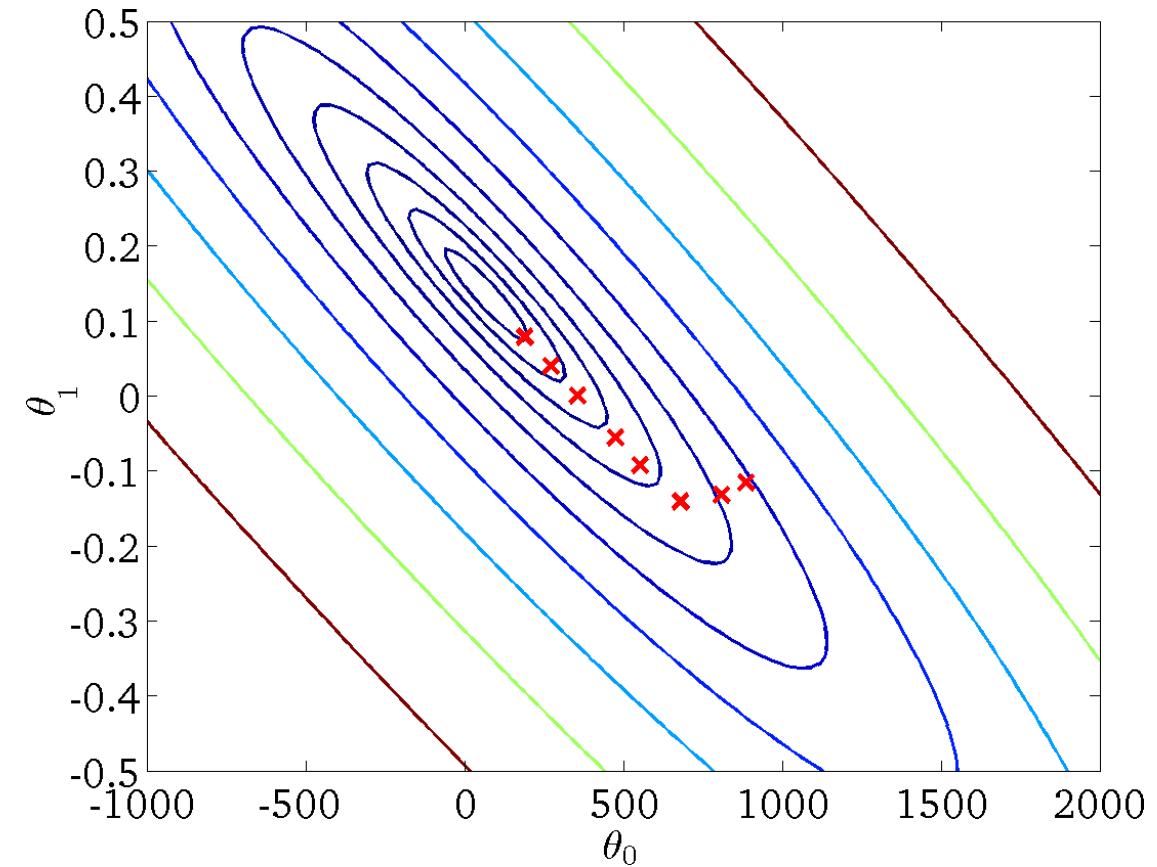
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



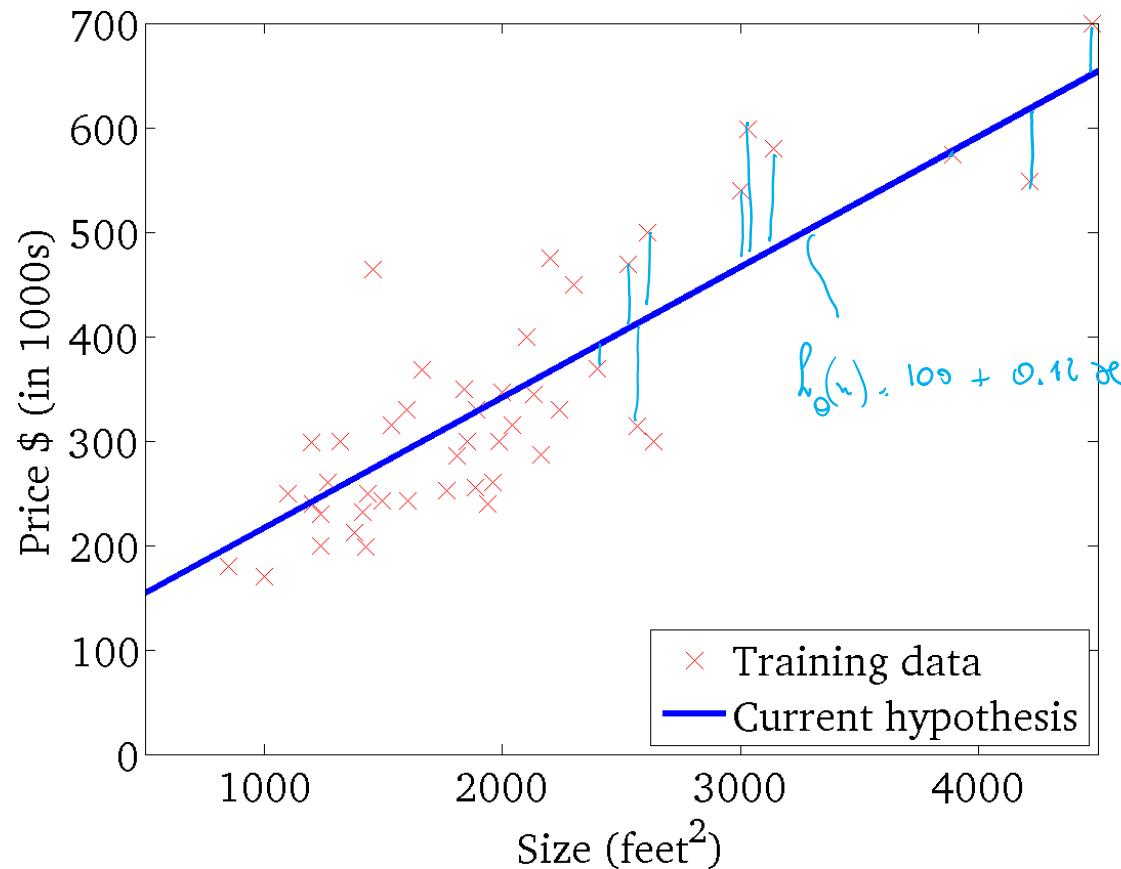
$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



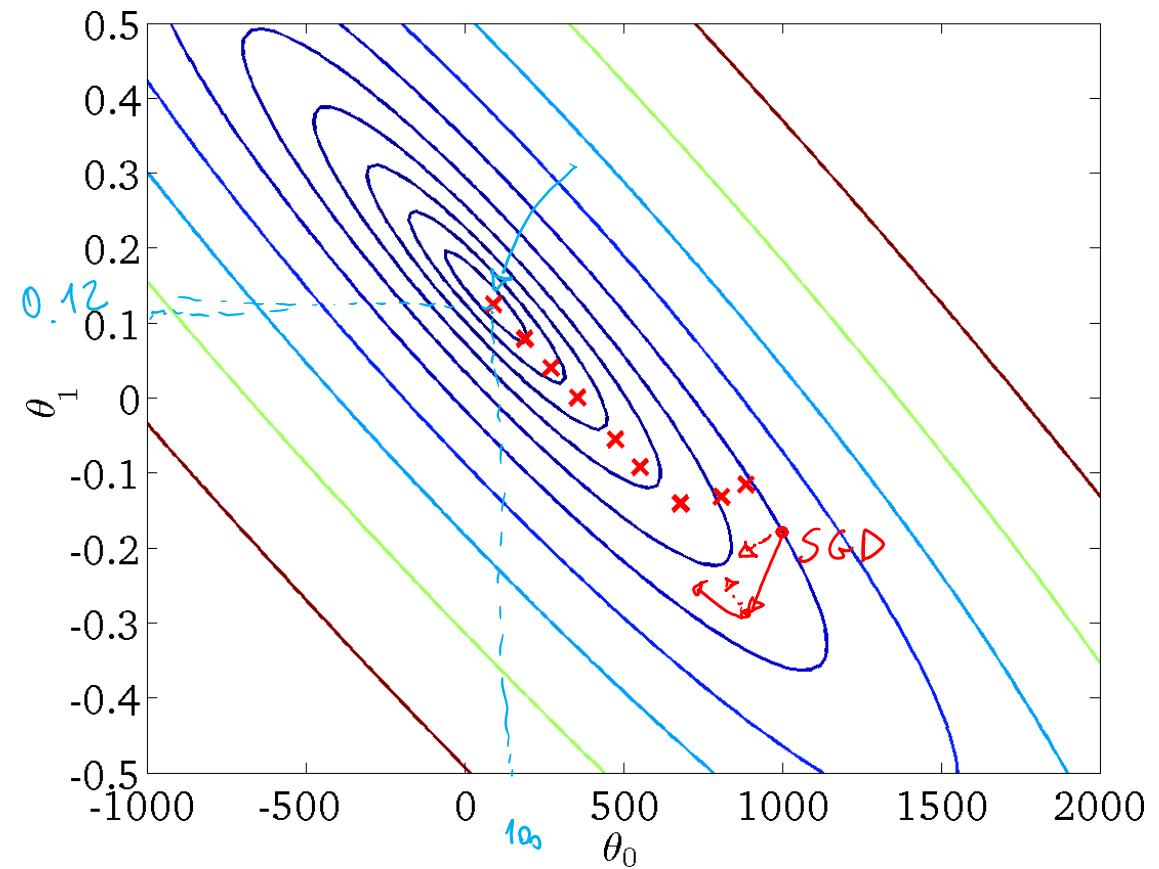
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



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

(function of the parameters θ_0, θ_1)



“Batch” Gradient Descent

“Batch”: Each step of gradient descent uses all the training examples.

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

Gradient Descent

- The gradient is the result of considering all dataset samples
- This may be very slow
- Can we do something more efficient?

Stochastic Gradient Descent (SGD)

- For large training sets, the evaluation the gradient over all samples may be expensive
- **Stochastic** or “**online**” gradient descent approximates the true gradient by a gradient at a single example
 - ▶ Pseudocode:
 - Choose an initial vector of parameters θ_0 and learning rate α
 - Repeat until convergence:
 - Randomly shuffle examples in the training set
 - For $k = 1, 2, \dots, m$, do:
 - ▶
$$\theta^{(i+1)} = \theta^{(i)} - \alpha \nabla J^{(k)}(\theta^{(i)})$$

$$\frac{\partial J}{\partial \theta} \approx \frac{1}{m} \sum_{i=1}^m (\dots)$$

Stochastic Gradient Descent (SGD)

- For large training sets, the evaluation the gradient over all samples may be expensive
- **Stochastic** or “**online**” gradient descent approximates the true gradient by a gradient at a single example
 - ▶ Pseudocode:
 - Choose an initial vector of parameters θ_0 and learning rate α
 - Repeat until convergence:
 - Randomly shuffle examples in the training set
 - For $k = 1, 2, \dots, m$, do:
 - ▶ $\theta^{(i+1)} = \theta^{(i)} - \alpha \nabla J^{(k)}(\theta^{(i)})$
 - Normally preferable: **mini-batch** gradient descent
 - ▶ Consider a mini-batch at each step
 - ▶ This normally results in smoother convergence
 - ▶ This normally is faster, thanks to vectorization libraries

*- SHUFFLE (PERMUTE ORDER)
- SELECT MINI-BATCHES
- ITERATE UP TO AN EPOCH
- REPEAT*

Gradient Descent: for Multiple Variables

Hypothesis: $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ Θ $= \theta^T x$

Cost function:

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$J(\underline{\Theta})$

Gradient descent:

Repeat {

$J(\underline{\Theta})$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$$

}

LEARNING RATE

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \underline{\Theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

n FEATURES

Gradient Descent

Previously (n=1):

Repeat {

$$\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})}_{\boxed{\frac{\partial}{\partial \theta_0} J(\theta)}}$$

$$\theta_1 := \theta_1 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}}_{\frac{\partial}{\partial \theta_1} J(\theta)}$$

(simultaneously update θ_0, θ_1)

}

New algorithm ($n \geq 1$):

Repeat {

$$\theta_j := \theta_j - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}}_{\frac{\partial}{\partial \theta_j} J(\theta)}$$

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

}

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \underbrace{x_0^{(i)}}_{\cancel{x_0^{(i)} = 1}}$$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

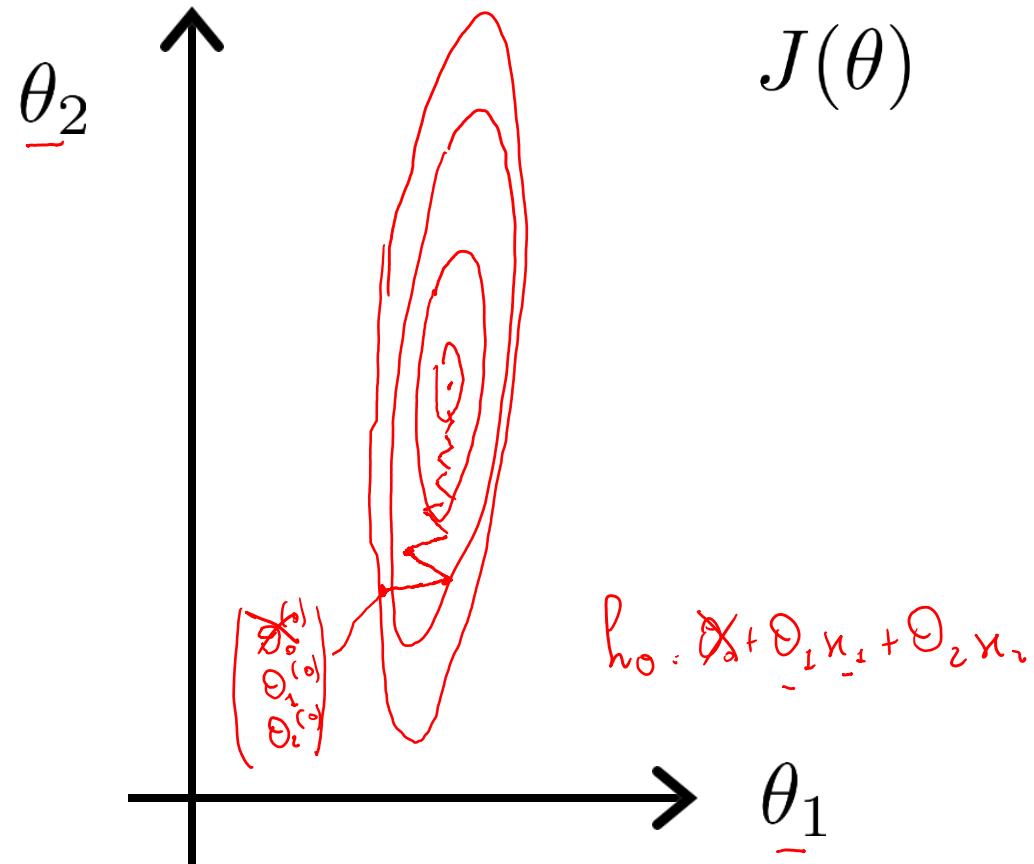
...

Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. \underline{x}_1 = size (0–2000 feet²)

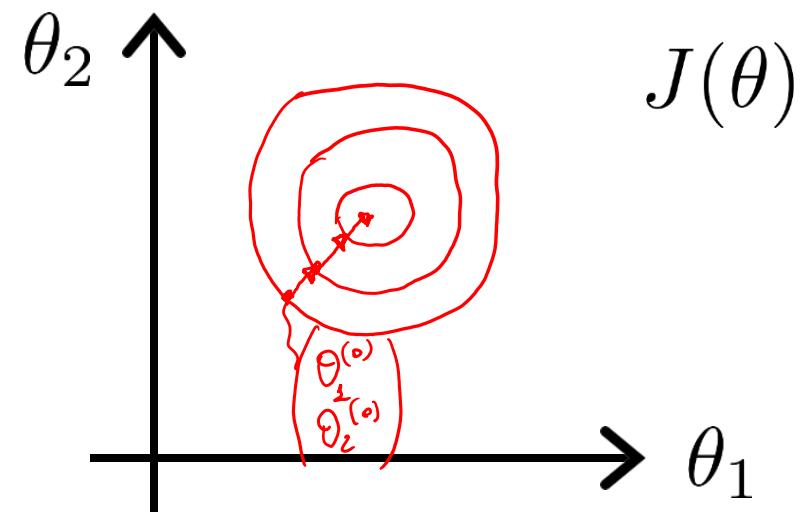
\underline{x}_2 = number of bedrooms (1–5)



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

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

$$-1 \leq x_j \leq 1$$



Feature Scaling

Get every feature into approximately a $-1 \leq x_i \leq 1$ range.

$$x_0 = 1$$

$$0 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-100 \leq x_3 \leq 500 \quad \times$$

$$-0.0001 \leq x_4 \leq 0.0006 \quad \times$$

Mean normalization

Replace x_i with $\frac{x_i - \mu_i}{\text{range}}$ to make features have approximately zero mean
(Do not apply to $x_0 = 1$).

E.g. $x_1 = \frac{\text{size} - 1090}{2000}$ $\mu_1 = 1090$ AVERAGE SIZE IN FEET²
RANGE

$$x_2 = \frac{\#\text{bedrooms} - 2}{5 - 4} \quad \mu_2 = 2 \quad 1 - 5$$

$$-0.5 \leq x_1 \leq 0.5, \quad -0.5 \leq x_2 \leq 0.5$$

$$x_j \leftarrow \frac{x_j - \mu_j}{s_j} \quad \begin{array}{l} \text{AVERAGE VALUE OF} \\ x_j \text{ IN TRAINING SET} \end{array} \quad j = 1 \dots n$$

RANGE ($\max - \min$)
(OR STANDARD DEVIATION)

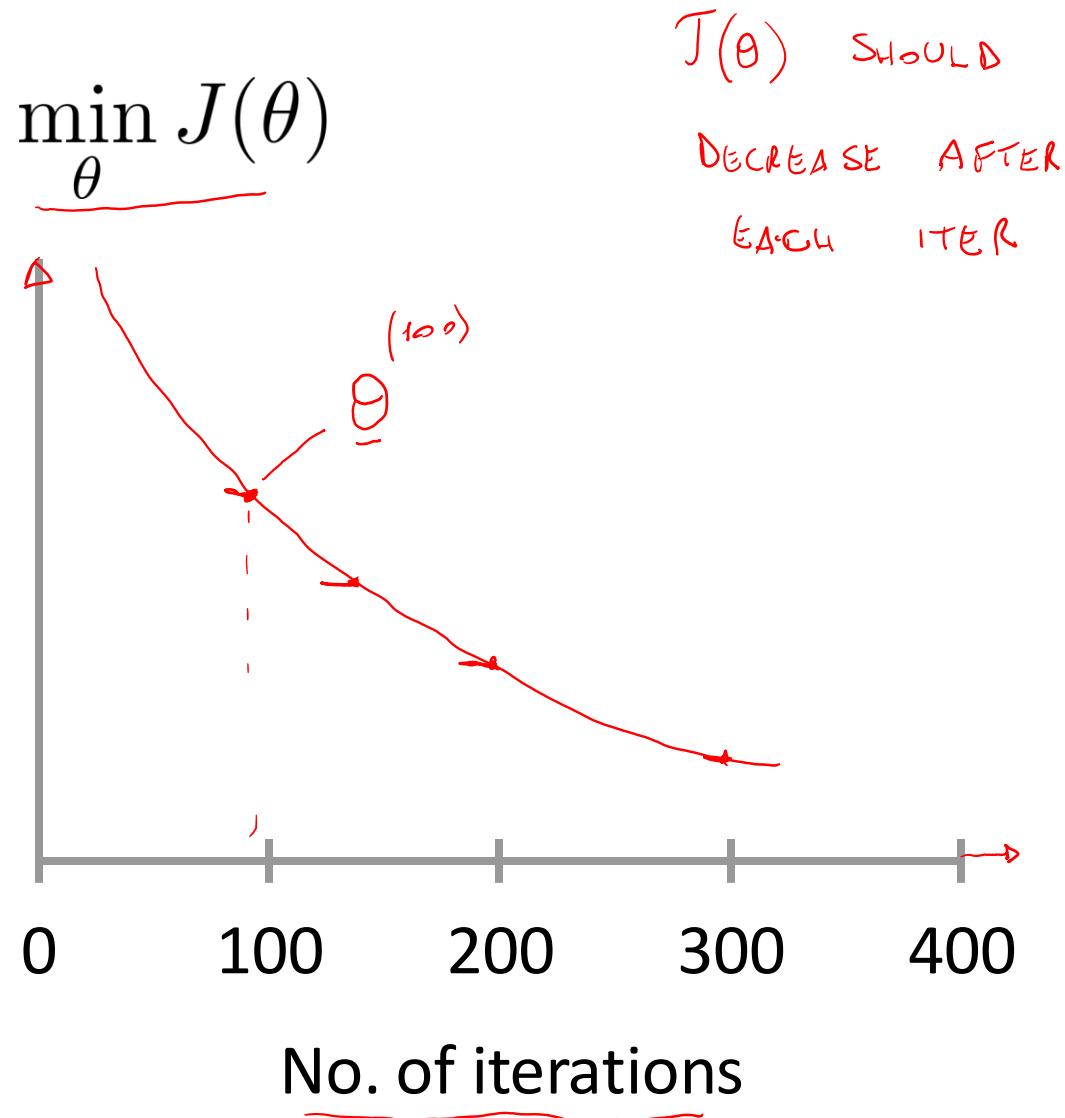
Gradient Descent: Learning Rate

Gradient descent

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

- “Debugging”: How to make sure gradient descent is working correctly.
- How to choose learning rate α .

Making sure gradient descent is working correctly.

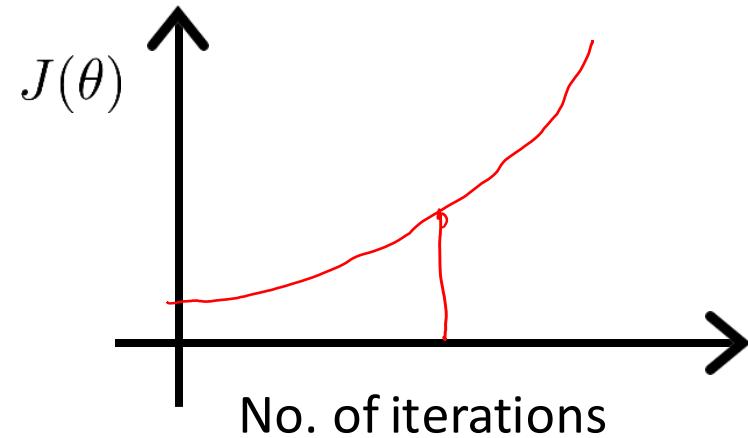


Example automatic convergence test:

Declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration.

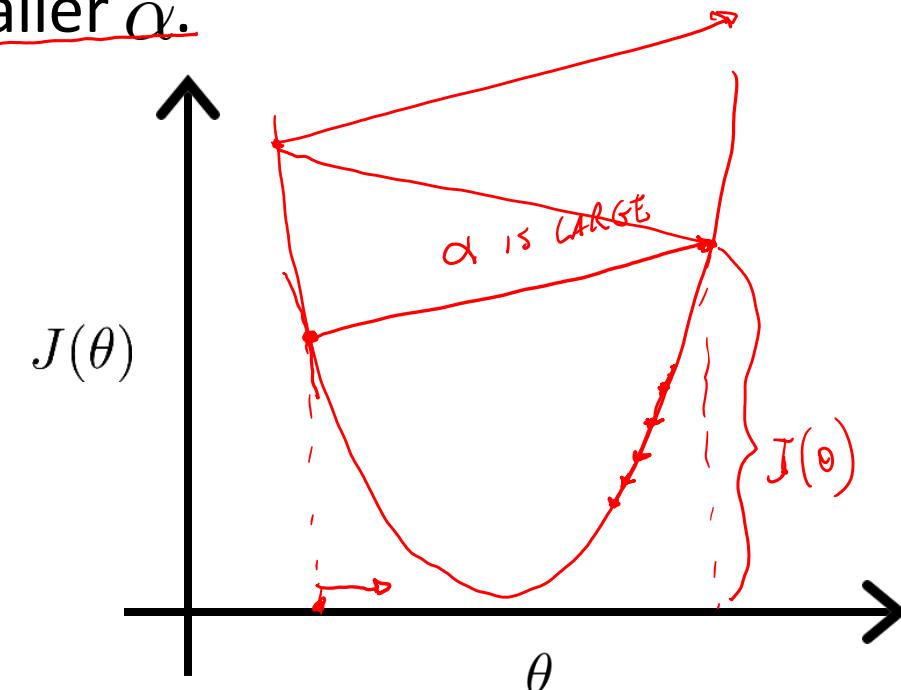
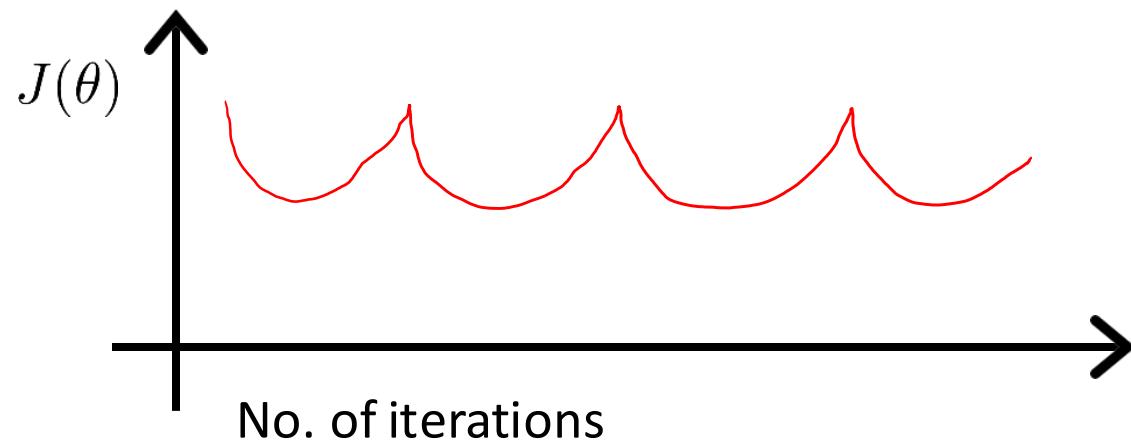
$$\epsilon$$

Making sure gradient descent is working correctly.



Gradient descent not working.

Use smaller α .



- For sufficiently small α , $J(\theta)$ should decrease on every iteration.
- But if α is too small, gradient descent can be slow to converge.

Summary:

- If α is too small: slow convergence.
- If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge.

To choose α , try

$$\dots, 0.001, \underbrace{0.003}_{3x}, 0.01, \underbrace{0.03}_{3x}, 0.1, \dots, 1, \dots$$

Linear Regression: Features and Polynomial regression

Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underline{\text{frontage}} + \theta_2 \times \underline{\text{depth}}$$

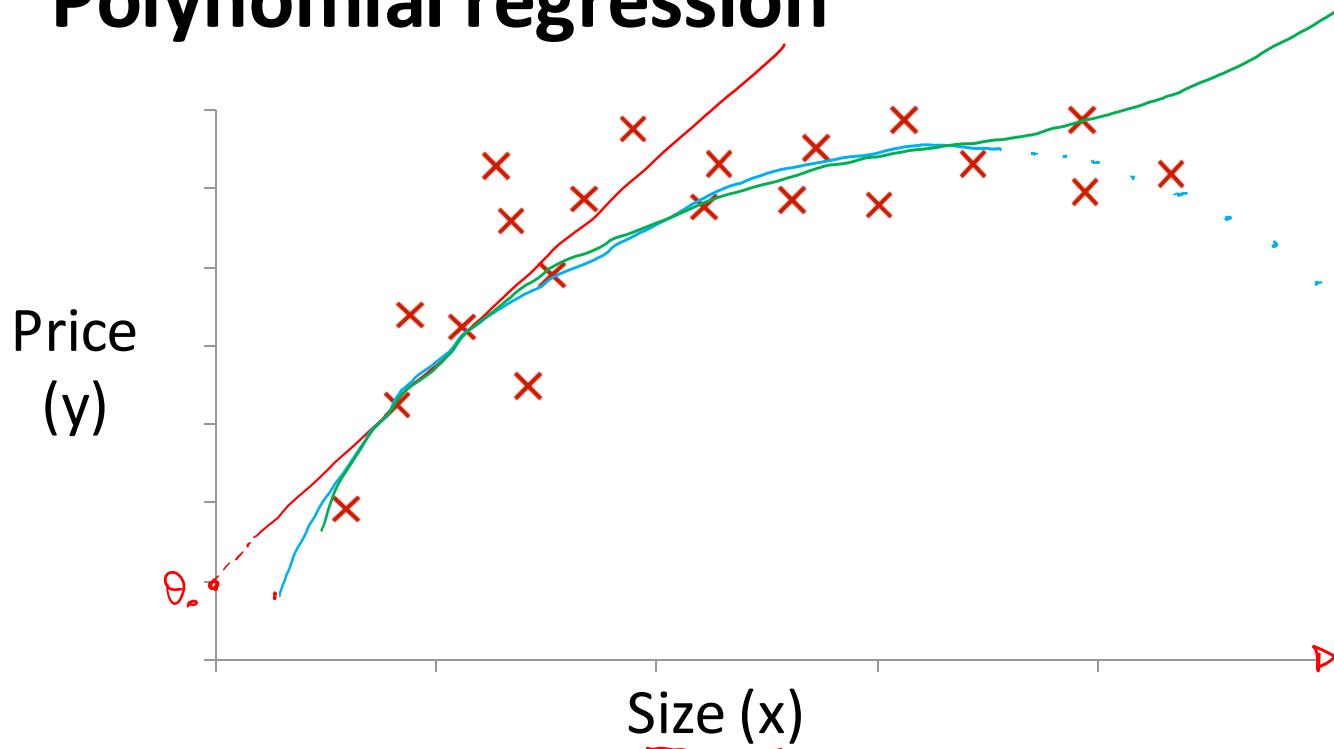
AREA

X , FRONTAGE X DEPTH

$$h_{\theta}(x) , \theta_0 + \theta_1 \times \nearrow \text{LAND AREA}$$



Polynomial regression



$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3 \end{aligned}$$

$$x_1 = (\text{size})$$

$$x_2 = (\text{size})^2$$

$$x_3 = (\text{size})^3$$

$$\theta_0 + \theta_1 x$$

$$\boxed{\theta_0 + \theta_1 x + \theta_2 x^2}$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \underline{\theta_3 x^3}$$

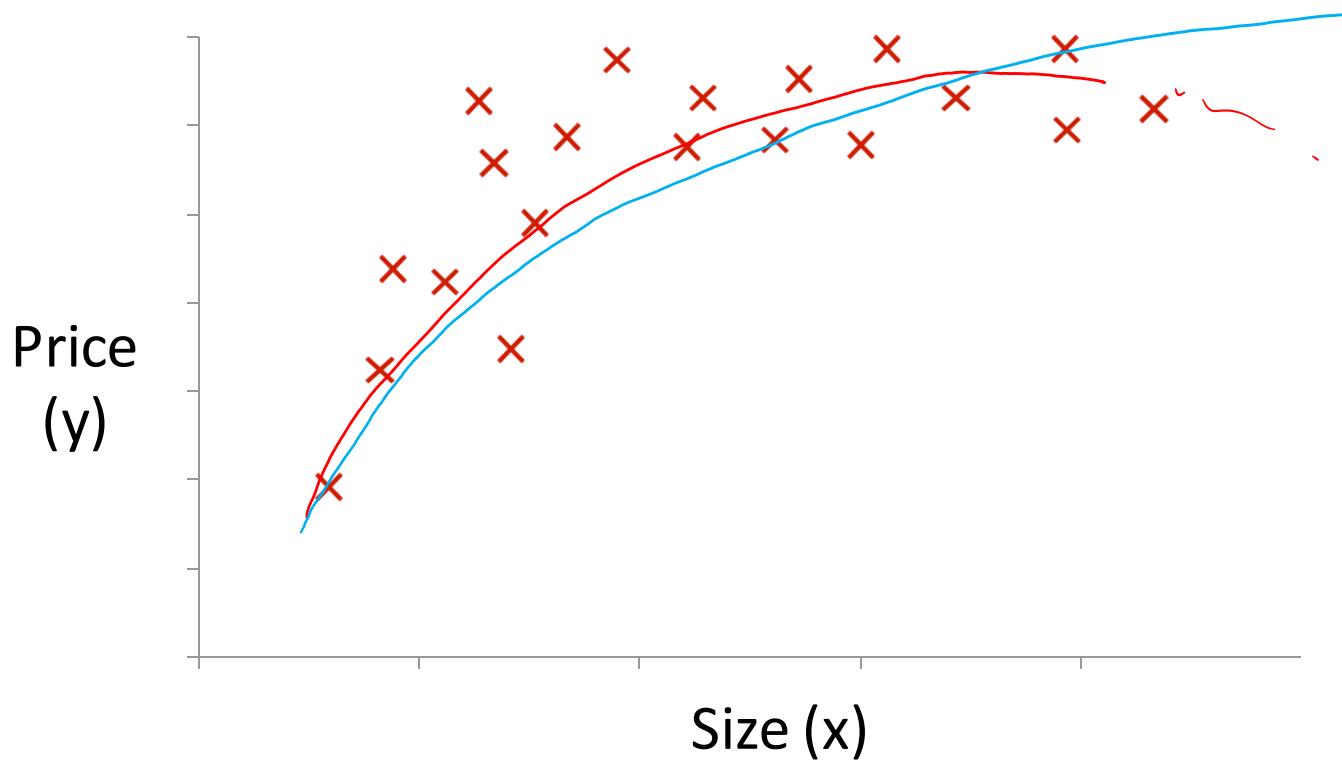
Size : 1 - 1000

Size² : 1 - 1,000,000

Size³ : 1 - 10⁹

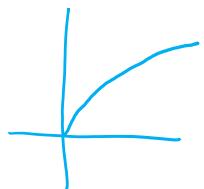
NEED TO SCALE
THE FEATURES

Choice of features



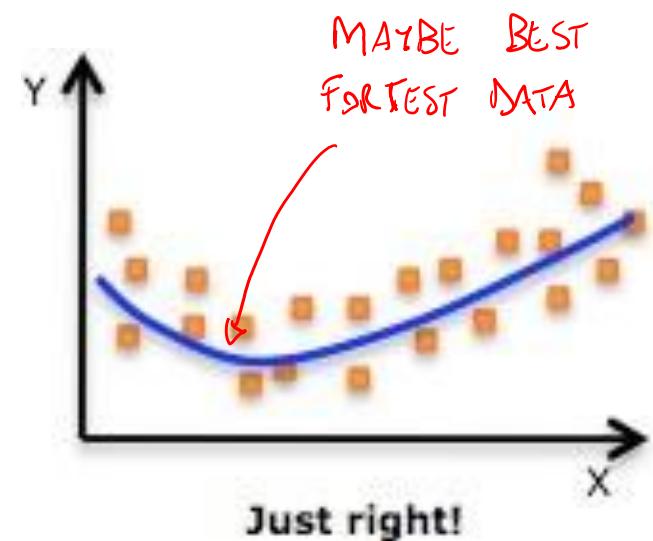
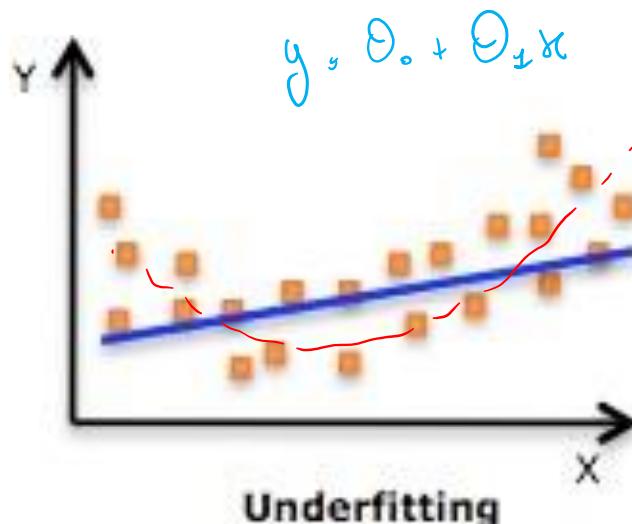
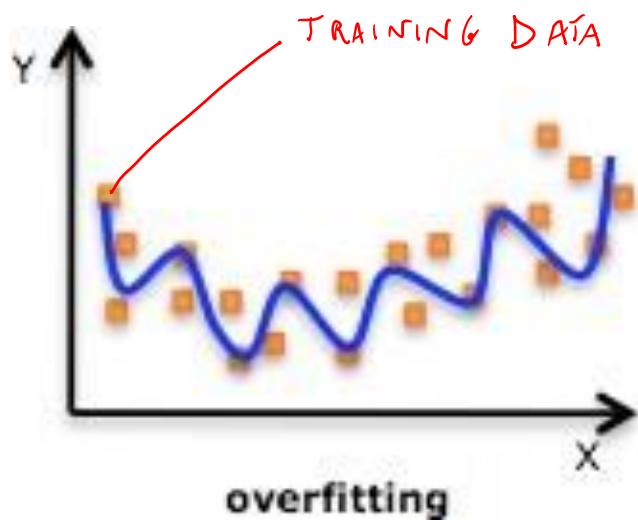
$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})}$$



Dangers of (Polynomial) Regression

Overfitting and Underfitting



Outline

- Linear regression
 - ▶ One or multiple variables
 - ▶ Cost function
 - Gradient descent, incl. batch/stochastic GD
 - Normal equation
 - MSE and Correlation
 - Locally-Weighted Regression
 - Probabilistic Interpretation of Least Squares
-

Normal Equation

$$\underline{\underline{\nabla}} f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}$$

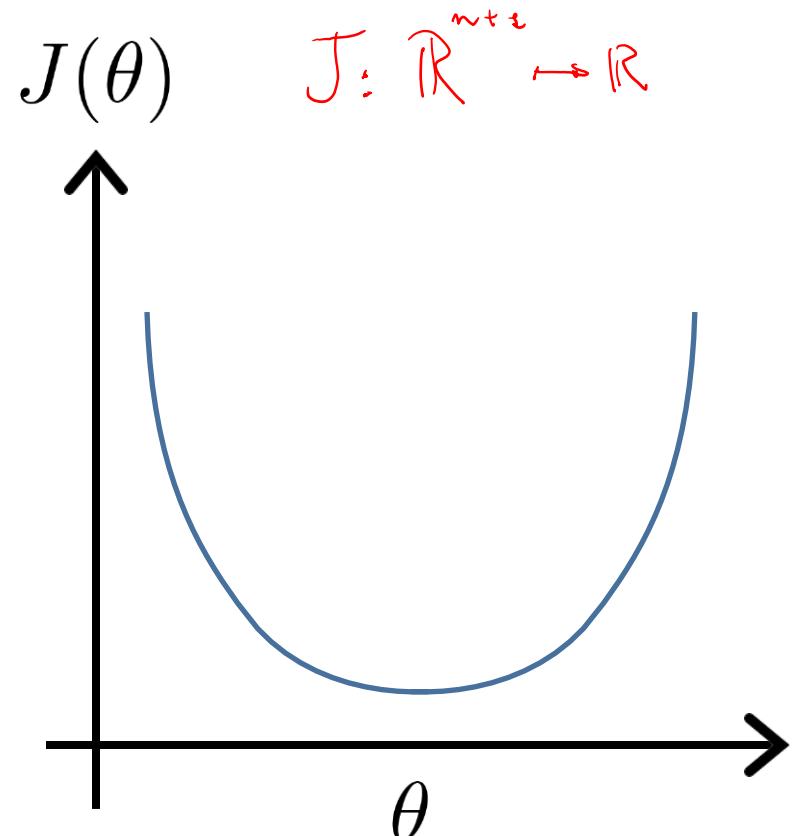
$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$$

Gradient Descent

$$\theta := \theta - \alpha \nabla_{\theta} J$$

$\underbrace{\quad}_{m \times n \text{ dim}}$ $\underbrace{\quad}_{m \times 1 \text{ dim}}$

$$\nabla_{\theta} J = \begin{bmatrix} \frac{\delta J}{\delta \theta_0} \\ \vdots \\ \frac{\delta J}{\delta \theta_n} \end{bmatrix} \in \mathbb{R}^{n+1}$$

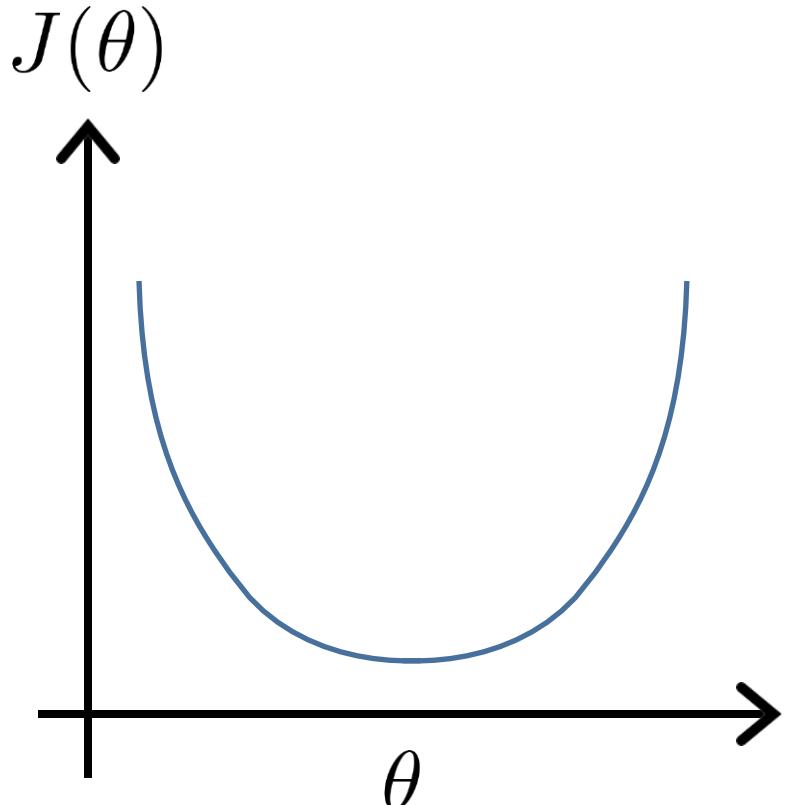


$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$$

Gradient Descent

$$\theta := \theta - \alpha \nabla_{\theta} J$$

$$\nabla_{\theta} J = \begin{bmatrix} \frac{\delta J}{\delta \theta_0} \\ \vdots \\ \frac{\delta J}{\delta \theta_n} \end{bmatrix} \in \mathbb{R}^{n+1}$$



Normal equation: what about solving for Θ analytically?

Useful notation

- For $f: \mathbb{R}^{m \times n} \mapsto \mathbb{R}$, define:

GRADIENT OF $f(A)$ w.r.t. A

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial}{\partial A_{11}} f & \cdots & \frac{\partial}{\partial A_{1n}} f \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial A_{m1}} f & \cdots & \frac{\partial}{\partial A_{mn}} f \end{bmatrix} \quad f(A) \in \mathbb{R}, A \in \mathbb{R}^{m \times n}$$

$\nabla_A f(A) \in \mathbb{R}^{m \times n}$

- Trace:

- If $A \in \mathbb{R}^{n \times n}$ $\text{tr } A = \sum_{i=1}^n A_{ii} \in \mathbb{R}$

Facts

- Some facts of matrix derivatives (without proof)

$$\text{tr } AB = \text{tr } BA$$

$$\text{tr } ABC = \text{tr } CAB = \text{tr } BCA$$

$$\underline{f(A) = \text{tr } BA} \in \mathbb{R} \quad \nabla_A \text{tr } AB = B^T$$

$$\text{tr } A = \text{tr } A^T$$

$$\text{If } a \in \mathbb{R} \quad \text{tr } a = a$$

$$\nabla_A \text{tr } ABA^T C = CAB + C^T A B^T$$

m examples $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; **n features.**

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

~~X~~ =

DESIGN MATRIX

$$\left[\begin{array}{c} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{array} \right]$$

→

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

$m \times (n+1)$

Cost function

$$X\Theta - y = \begin{bmatrix} h(x^{(1)}) - y^{(1)} \\ \vdots \\ h(x^{(m)}) - y^{(m)} \end{bmatrix} \quad \text{where} \quad \Theta^T \underline{x^{(i)}} = \sum_{j=0}^n \Theta_j x_j^{(i)}$$

DESIGN
MATRIX

Recall $\underline{z^T z} = \sum_i z_i^2$

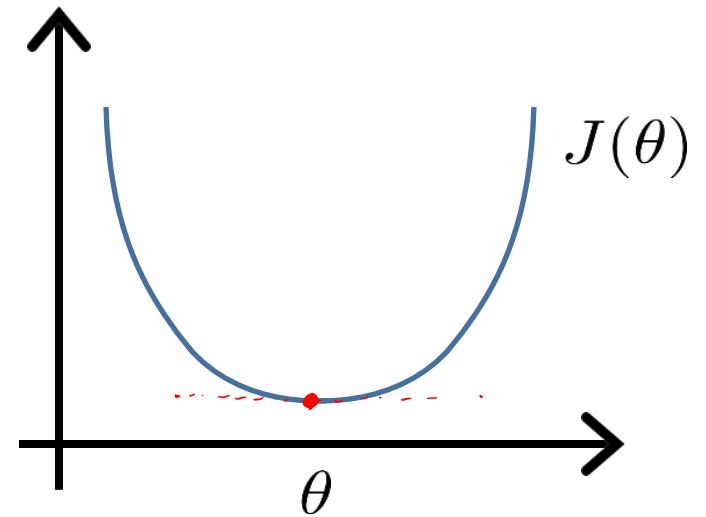
$$\underline{\frac{1}{2}(X\Theta - y)^T(X\Theta - y)} = \frac{1}{2} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 = \underline{J(\Theta)} \quad \theta \in \mathbb{R}^{n+1}$$

Intuition: If 1D ($\theta \in \mathbb{R}$)

$$J(\theta) = a\theta^2 + b\theta + c$$

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

SOLVE FOR θ



Solve for Θ analytically

$$\theta^* = \arg \min_{\theta} J(\theta) \quad \text{attained when} \quad \nabla_{\theta} J(\theta) = 0$$

Expanding $J(\theta)$:

$$\begin{aligned} & \nabla_{\theta} \frac{1}{2} (X\theta - y)^T (X\theta - y) \quad \text{tr } a \cdot a \quad a^T \cdot a \\ &= \frac{1}{2} \nabla_{\theta} \text{tr} \left(\theta^T X^T X \theta - \underbrace{\theta^T X^T y}_{\text{SCALAR}} - y^T X \theta + \underbrace{y^T y}_{\text{SCALAR}} \right) \\ &= \frac{1}{2} \left[\nabla_{\theta} \text{tr} \theta \theta^T X^T X - \nabla_{\theta} \text{tr} y^T X \theta - \nabla_{\theta} \text{tr} y^T X \theta \right] \end{aligned}$$

Solve for Θ analytically

$$\frac{1}{2} \left[\nabla_{\Theta} \text{tr} \Theta \Theta^T X^T X - \nabla_{\Theta} \text{tr} y^T X \Theta - \nabla_{\Theta} \text{tr} y^T X \Theta \right]$$

Recall $\nabla_A \text{tr} ABA^T C = CAB + C^T AB^T$ $\nabla_A \text{tr} AB = B^T$ $\text{tr} AB = \text{tr} BA$

$$\nabla_{\Theta} \text{tr} \underbrace{\Theta \text{I}}_A \underbrace{\Theta^T}_A \underbrace{X^T}_B \underbrace{X}_C = \underbrace{X^T X}_C \underbrace{\Theta \text{I}}_A + \underbrace{X^T X}_C \underbrace{\Theta \text{I}}_B$$

$$\nabla_{\Theta} \text{tr} \underbrace{y^T}_A \underbrace{X \Theta}_B = X^T y$$

Solve for Θ analytically

$$\begin{aligned}\underline{\nabla_{\theta} J} &= \frac{1}{2} \left(X^T X \theta + X^T X \theta - 2 X^T \bar{y} \right) \\ &= \underline{X^T X \theta} - X^T \bar{y} = 0\end{aligned}$$

$$\underbrace{X^T X \theta}_{=} = X^T \bar{y} \quad \text{Normal equation}$$

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

Examples: $m = 4$.

x_0	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
	x_1	x_2	x_3	x_4	y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \in \mathbb{R}^{m \times (m+1)}$

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix} \in \mathbb{R}^m$

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

m training examples, n features.

Gradient Descent

- Need to choose α .
- Needs many iterations.
- Works well even when n is large.

$$n = 10^6$$

Normal Equation

- No need to choose α .
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$
- Slow if n is very large.

$$n = 100 \dots n = 10,000$$

Normal equation

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

- What if $X^T X$ is non--invertible? (singular/ degenerate)

What if $X^T X$ is non--invertible?

- Redundant features (linearly dependent).

E.g. $x_1 = \text{size in feet}^2$

1 m. 3.28 FEET

~~$x_2 = \text{size in m}^2$~~

$$x_1 = (3.28)^2 x_2$$

- Too many features (e.g. $m \leq n$).

- Delete some features, or use regularization.

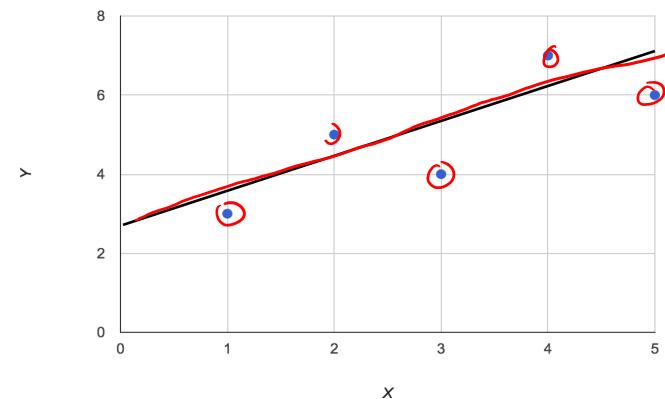
Outline

- Linear regression
 - ▶ One or multiple variables
 - ▶ Cost function
 - Gradient descent, incl. batch/stochastic GD
 - Normal equation
 - MSE and Correlation
 - Locally-Weighted Regression
 - Probabilistic Interpretation of Least Squares
-

Linear Regression and Correlation

Summary So Far

- Given: Set of known (x,y) points
 - Find: function $f(x)=ax+b$ that “best fits” the known points, i.e., $f(x)$ is close to y
 - Use function to predict y values for new x 's
- Also can be used to test correlation



Correlation and Causation

Correlation - Values track each other

- Height and Shoe Size
- Grades and Entrance Exam Scores

Causation - One value directly influences another

- Education Level → Starting Salary
- Temperature → Cold Drink Sales

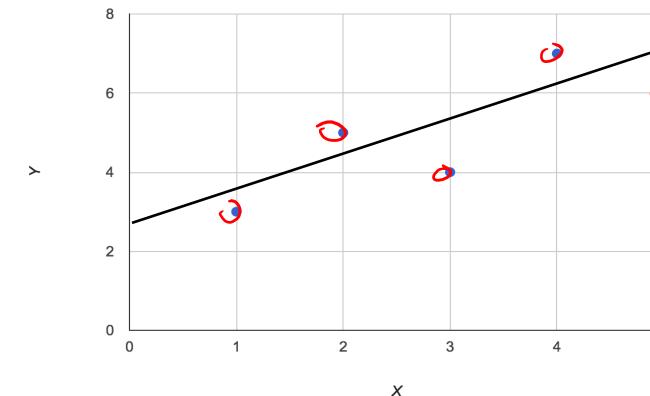
Correlation and Causation (from Overview)

Correlation - Values track each other

- Height and Shoe Size
- Grades and Entrance Exam Scores

Find: function $f(x)=ax+b$ that “best fits” the known points, i.e., $f(x)$ is close to y

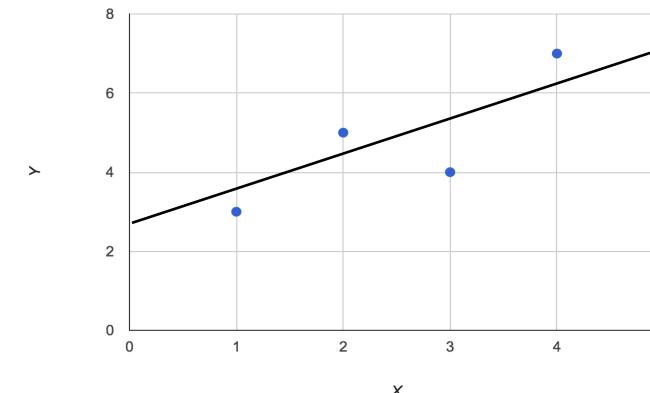
The better the function fits the points, the more correlated x and y are



Regression and Correlation

The better the function fits the points,
the more correlated x and y are

- Linear functions only
- Correlation - Values track each other
 - Positively - when one goes up the other goes up
- Also negative correlation
 - When one goes up the other goes down
 - Latitude versus temperature
 - Car weight versus gas mileage
 - Class absences versus final grade



Calculating Simple Linear Regression

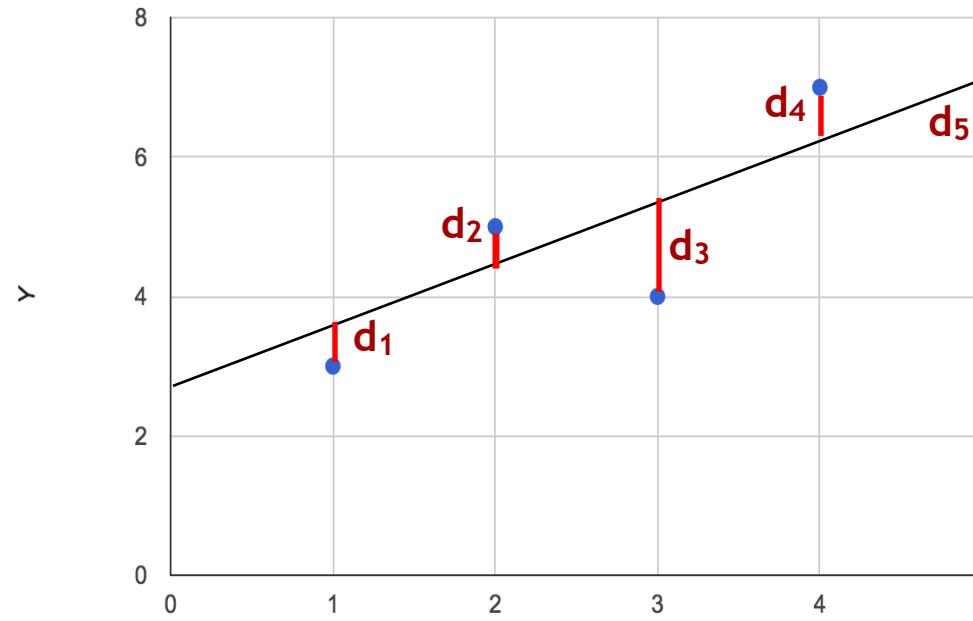
Method of least squares

- Given a point and a line, the error for the point is its vertical distance d from the line, and the squared error is d^2
- Given a set of points and a line, the sum of squared error (SSE) is the sum of the squared errors for all the points
- Goal: Given a set of points, find the line that minimizes the SSE

J

Calculating Simple Linear Regression

Method of least squares

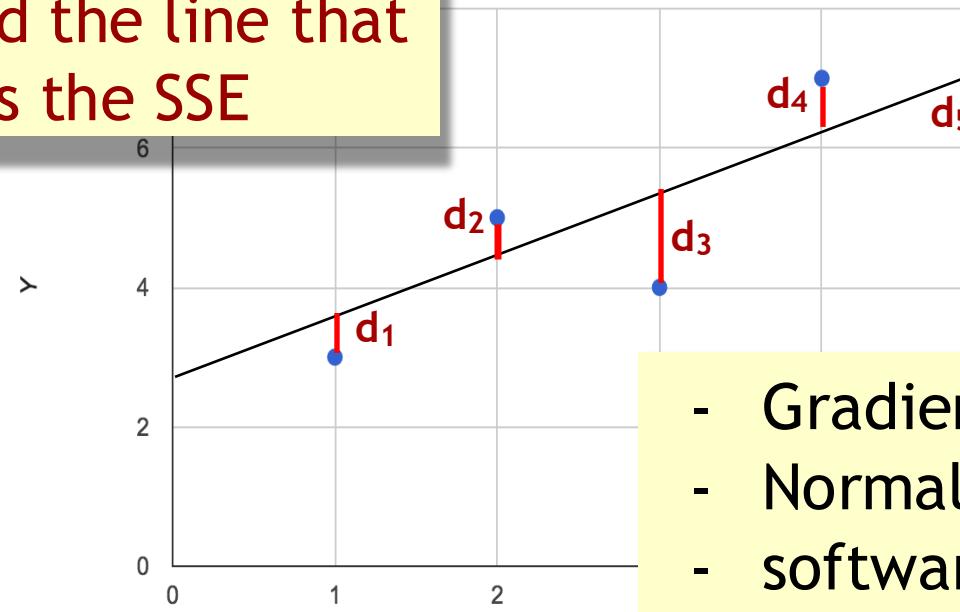


$$SSE = d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2$$

Calculating Simple Linear Regression

Method of least squares

Goal: Find the line that minimizes the SSE



- Gradient Descent
- Normal equation
- software packages,
e.g. Numpy polyfit

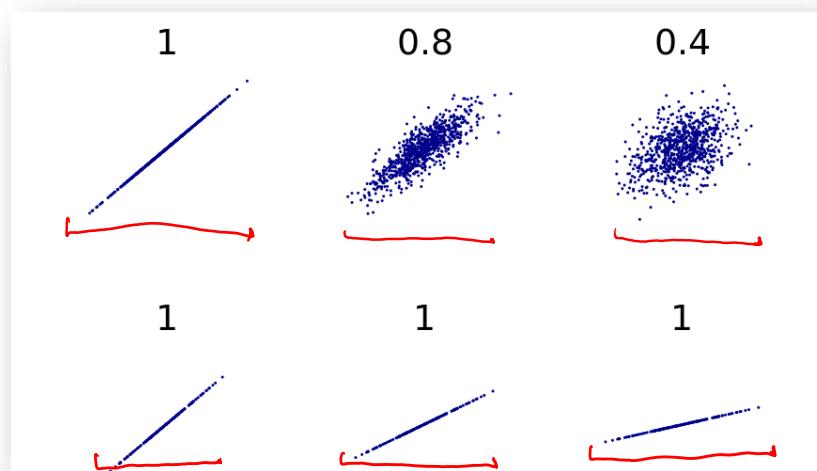
$$SSE = d_1^2 + d_2^2 + d_3^2 + d_4^2 + d_5^2$$

Measuring Correlation

More help from software packages...

Pearson's Product Moment Correlation (PPMC)

- “Pearson coefficient”, “correlation coefficient”
- Value r between 1 and -1
 - 1 maximum positive correlation
 - 0 no correlation
 - 1 maximum negative correlation



Swapping x and y axes
yields same values

Measuring Correlation

More help from software packages...

Pearson's Product Moment Correlation (PPMC)

- “Pearson coefficient”, “correlation coefficient”
- Value r between 1 and -1
 - 1 maximum positive correlation
 - 0 no correlation
 - 1 maximum negative correlation

Coefficient of determination

- r^2 , R^2 , “R squared”
- Measures fit of any line/curve to set of points
- Usually between 0 and 1
- For simple linear regression $R^2 = \text{Pearson}^2$

“The better the function fits the points, the more correlated x and y are”

Correlation Game

<http://aionet.eu/corguess> (*)

Try to get:

Right answers ≥ 10 , Guesses \leq Right answers $\times 2$

Anti-cheating: Pictures = Right answers + 1

(*) Improved version of “Wilderdom correlation guessing game” thanks to Poland participant Marcin Piotrowski

Other correlation games:

<http://guessthecorrelation.com/>

<http://www.rossmanchance.com/applets/GuessCorrelation.html>

<http://www.istics.net/Correlations/>

Outline

- Linear regression
 - ▶ One or multiple variables
 - ▶ Cost function
 - Gradient descent, incl. batch/stochastic GD
 - Normal equation
 - MSE and Correlation
 - Locally-Weighted Regression
 - Probabilistic Interpretation of Least Squares
-

Locally-Weighted Regression

Recap

m **examples** $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; n **features**.

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad y^{(i)} \in \mathbb{R} \quad x_0 = 1$$

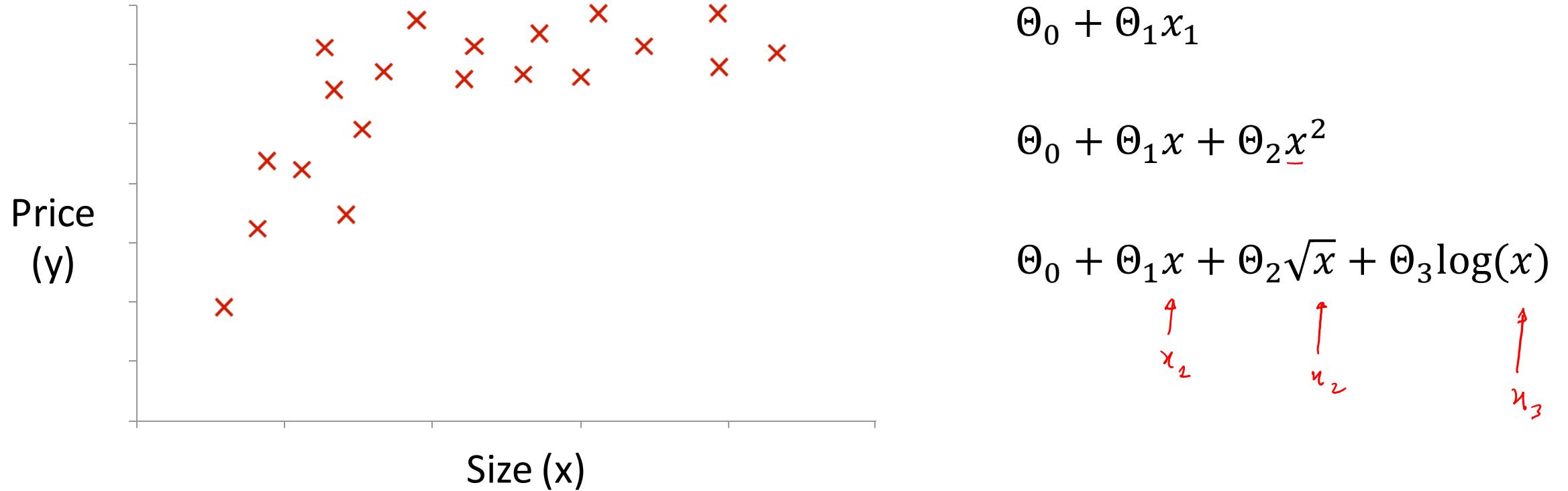
$x^{(i)}$ FEATURES FOR
 i -th TRAINING
EXAMPLE

$h_{\Theta}(x^{(i)})$ PREDICTION FOR
 i -th TRAINING
EXAMPLE

x NEW EXAMPLE
(TEST)

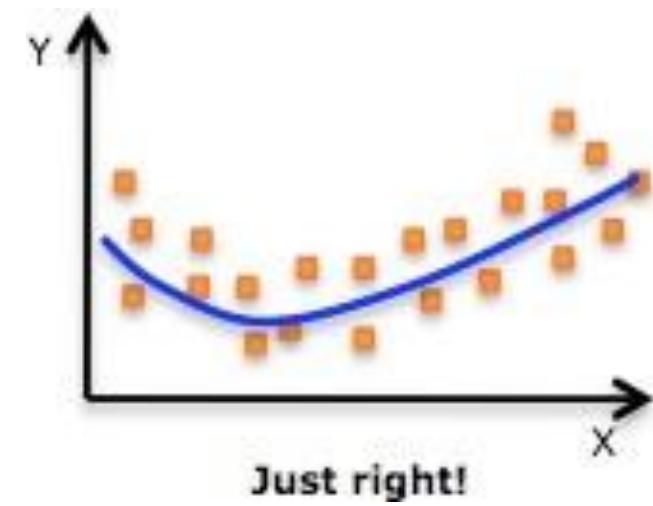
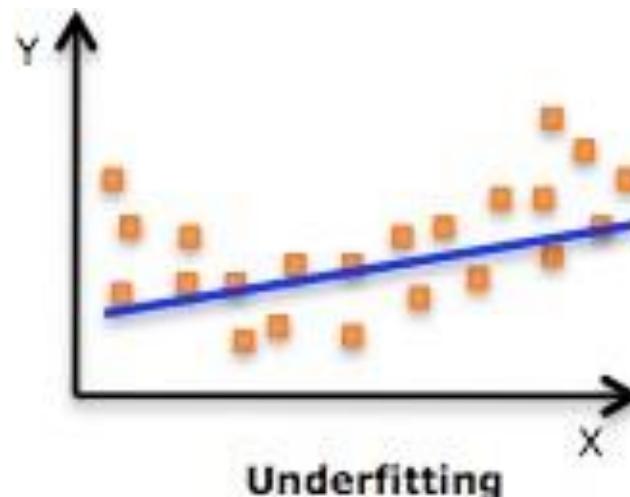
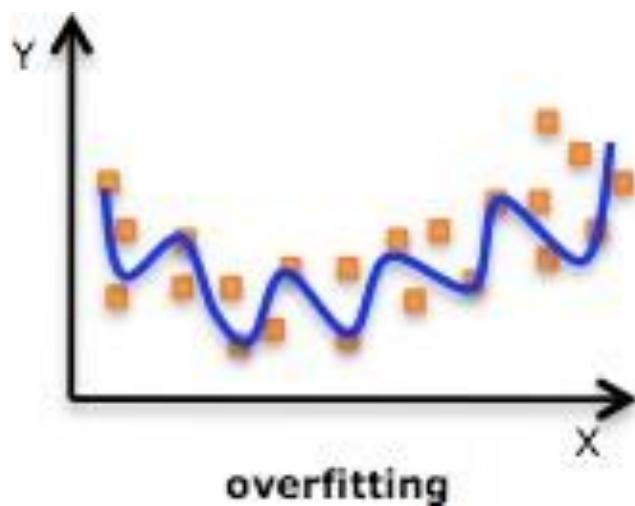
$$\underline{h_{\Theta}(x)} = \sum_{j=0}^n \Theta_j x_j = \Theta^T x \quad J(\Theta) = \frac{1}{2} \sum_{i=1}^m (h_{\Theta}(x^{(i)}) - y^{(i)})^2$$

Recap: choice of features



Recap: dangers of polynomial regression

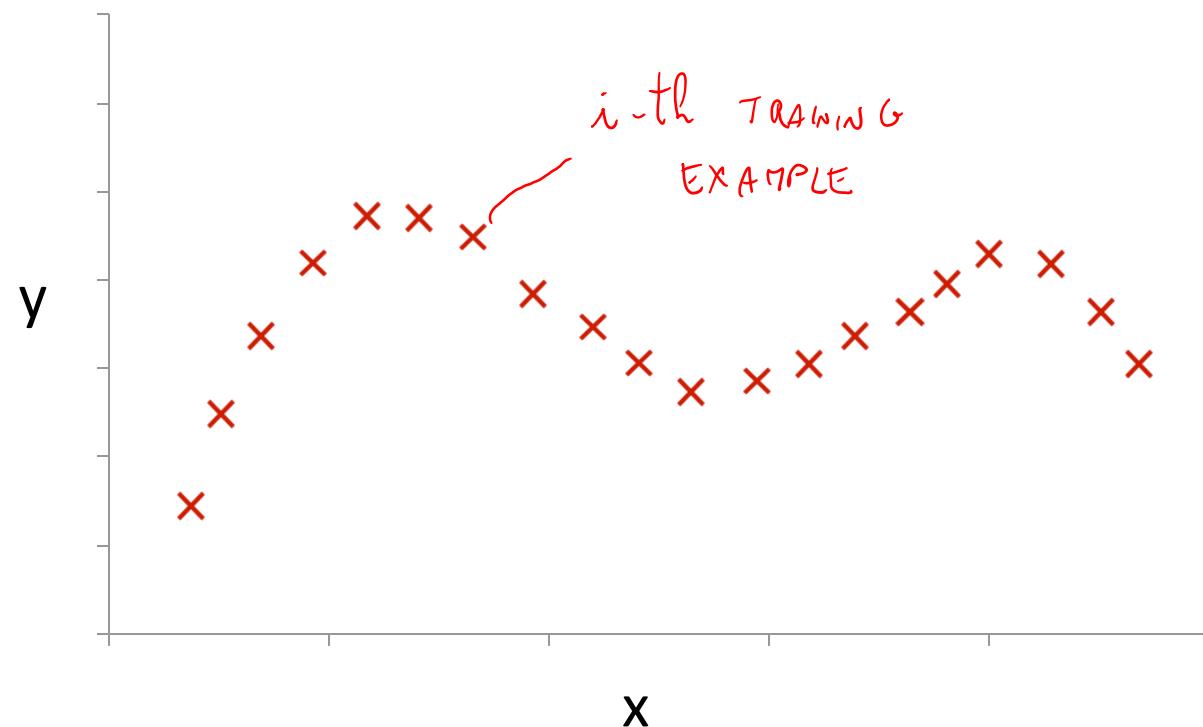
Overfitting and Underfitting



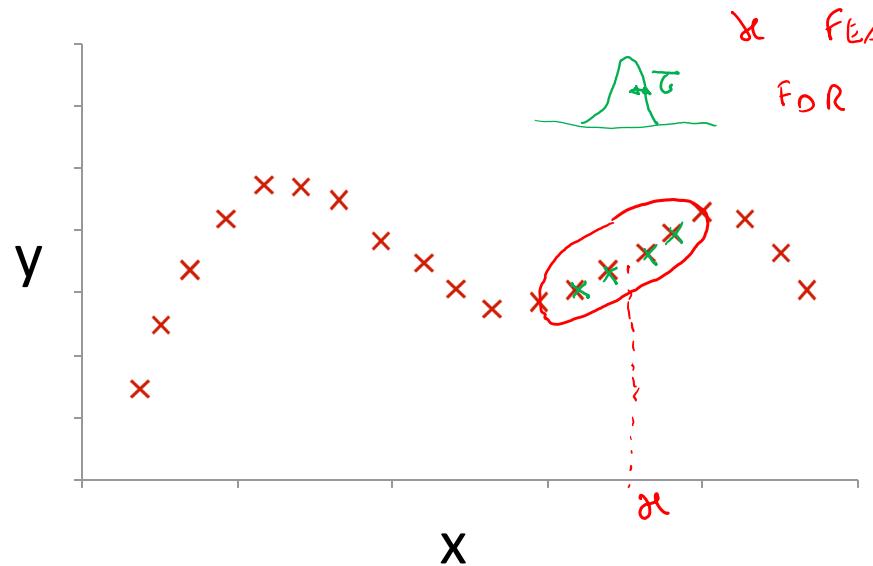
Locally-weighted regression

- “Parametric learning algorithm”: fixed set of parameters
- “Non-parametric learning algorithm”: parameters grow with the data (more data to keep in memory)
- Locally-weighted regression
 - ▶ also named Loess or Lowess

$$LR: \hat{y} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



Locally-weighted regression



LR:

TO EVALUATE h AT CERTAIN x

$$\text{FIT } \theta \text{ TO MINIMIZE} \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$$

RETURN $\theta^T x$

LOCALLY WEIGHTED REGRESSION

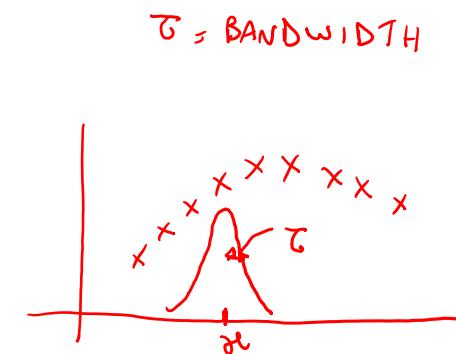
FIT θ TO MINIMIZE

$$\sum_{i=1}^m w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

$$w^{(i)} = \exp\left(-\frac{\|x^{(i)} - x\|^2}{2\sigma^2}\right)$$

IF $\|x^{(i)} - x\|$ SMALL $w^{(i)} \approx 1$

$\|x^{(i)} - x\|$ LARGE $w^{(i)} \approx 0$



Outline

- Linear regression
 - ▶ One or multiple variables
 - ▶ Cost function
 - Gradient descent, incl. batch/stochastic GD
 - Normal equation
 - MSE and Correlation
 - Locally-Weighted Regression
 - Probabilistic Interpretation of Least Squares
-

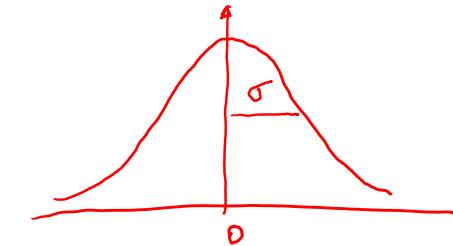
Probabilistic Interpretation of Least Squares

Why Least Squares

- Assume $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$ ← "error": UNMODELED EFFECTS, RANDOM NOISE

$$\epsilon^{(i)} \sim N(0, \sigma^2)$$

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$



ASSUMPTION: $\epsilon^{(i)}$ ARE IID

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$



PARAMETERIZED BY

$$y^{(i)} | x^{(i)}; \theta \sim N(\theta^T x^{(i)}, \sigma^2)$$

Why Least Squares

Contd

- Likelihood $L(\theta) = P(\vec{y}|X; \theta)$

$$\begin{aligned} L(\theta) &= P(\vec{y} | X; \theta) \\ &= \prod_{i=1}^m P(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

LOG LIKELIHOOD

$$\begin{aligned} \ell(\theta) &= \log L(\theta) \\ &= \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp(\dots) \end{aligned}$$

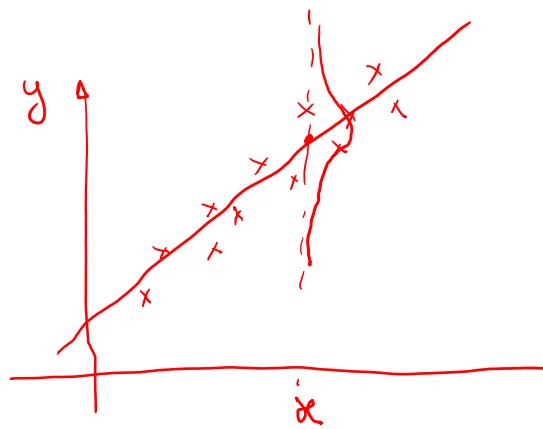
$$\begin{aligned} &= \sum_{i=1}^m \left[\log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp(\dots) \right] \\ &= m \log \frac{1}{\sqrt{2\pi}\sigma} + \sum_{i=1}^m -\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2} \end{aligned}$$

MLE: MAXIMUM LIKELIHOOD ESTIMATION
CHOOSE θ TO MAXIMIZE $L(\theta)$

MINIMIZE $\frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 = J(\theta)$



More on the Probabilistic Interpretation



$$p(y | \alpha, \theta) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(y - \theta^\top \alpha)^2}{2\sigma^2}}$$

References

- Chapter 3.1 in [Bishop, 2006. Pattern Recognition and Machine Learning]
- Chapter 9 in [Deisenroth Faisal Ong Book 2020 Mathematics for Machine Learning]

Thank you

Acknowledges: slides and material from Andrew Ng, Eric Xing, Matthew R. Gormley, Jessica Wu

