

Machine Learning



BITS Pilani
Pilani Campus

Dr. Monali Mavani



Machine Learning



Disclaimer and Acknowledgement

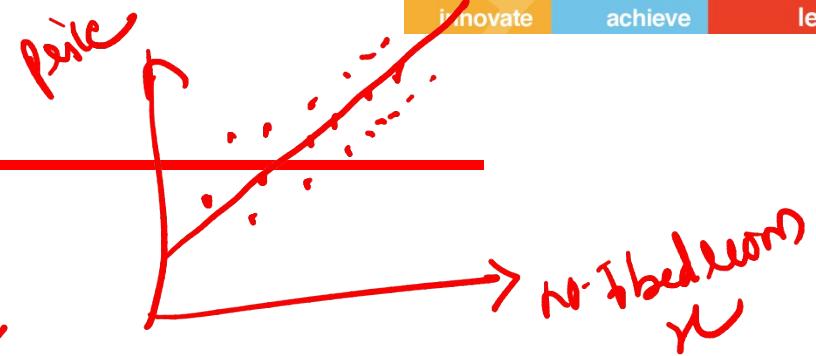


- The content for these slides has been obtained from books and various other source on the Internet
- I here by acknowledge all the contributors for their material and inputs.
- I have provided source information wherever necessary
- I have added and modified the content to suit the requirements of the course

Session Content

Linear models for regression

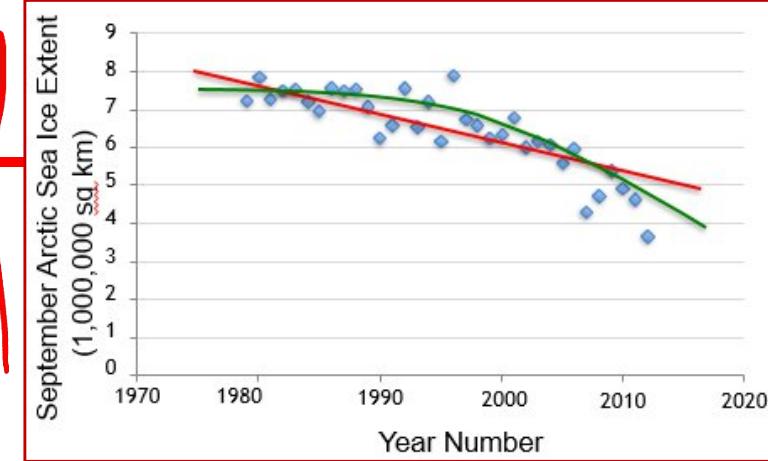
- Linear Regression
- Linear Regression Cost Function
- Learning Model Parameters - Gradient Descent → Optimization
- Learning Model Parameters – Closed Solution
- Linear Basis Function Models
- Quality of fit
- Evaluation metrics
-



Analytical

Linear Regression

d -dimensional features
 $d = n - 1$



Given:

- Data $\underline{X} = \{\underline{x}^{(1)}, \dots, \underline{x}^{(n)}\}$ where $\underline{x}^{(i)} \in \mathbb{R}^d$
 - Corresponding labels $\underline{y} = \{y^{(1)}, \dots, y^{(n)}\}$ where $y^{(i)} \in \mathbb{R}$
- $y = f(x) \rightarrow \text{true function}$
- $x_1 \ x_2$

Wish to learn $f : \underline{X} \rightarrow \underline{Y}$,

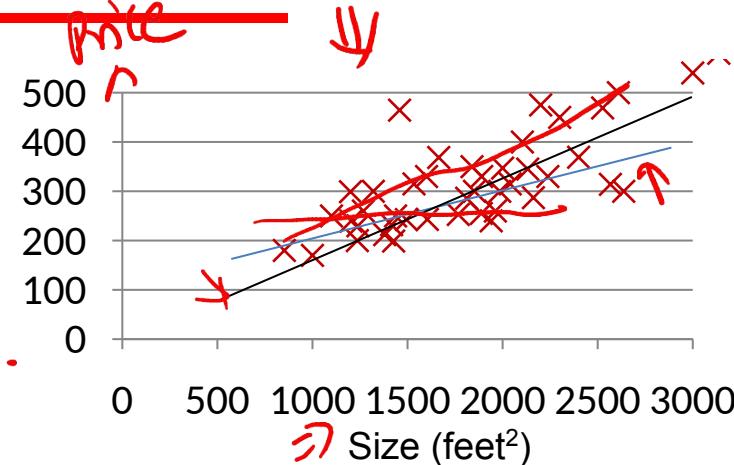
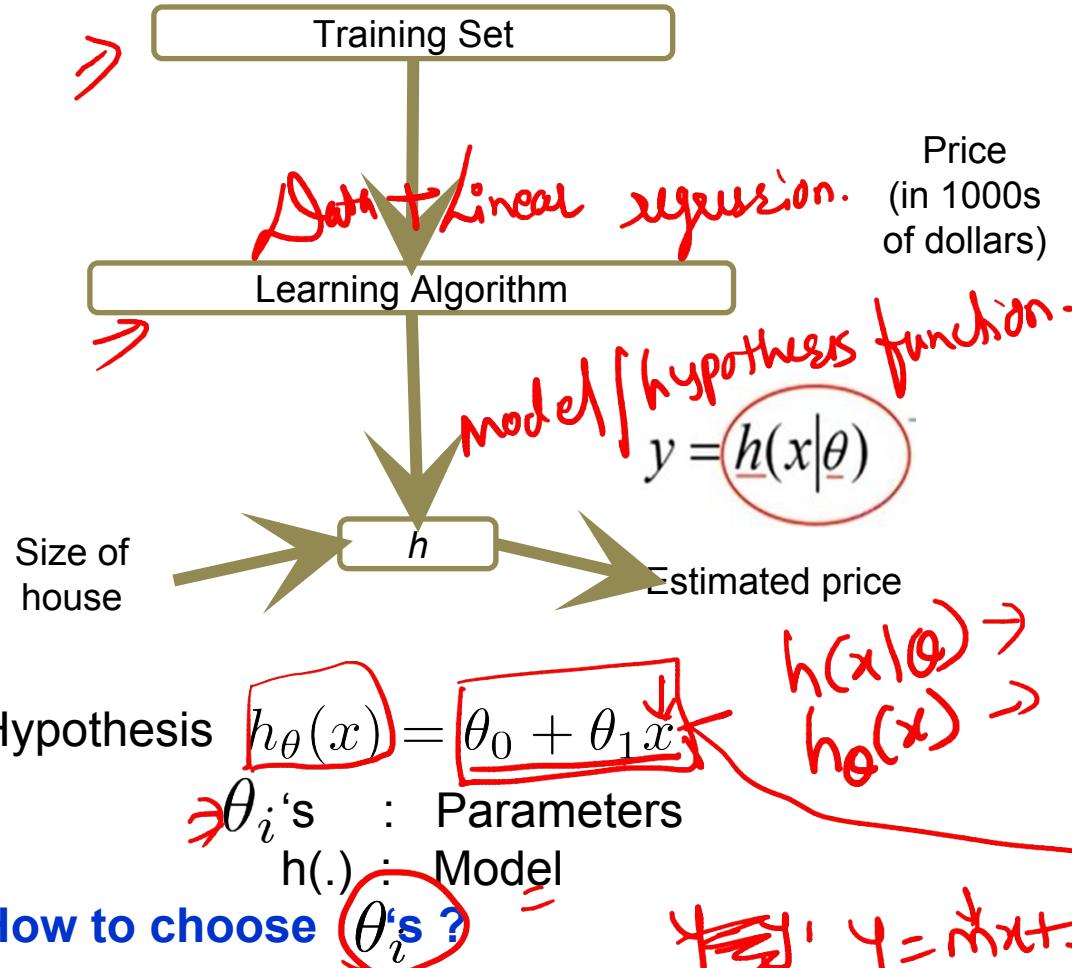
For each input observation $x^{(i)}$, a vector of features $[x_1, x_2, \dots, x_n]$.

j^{th} feature of i^{th} training example: $x_j^{(i)}$

\rightarrow training example
 \rightarrow feature

Machine Learning Process : Interpretation

$y = f(x) \Rightarrow y = h(x|\theta) \rightarrow$ general symbol for parameters



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

$y = mx + b \Rightarrow y = \theta_0 x + \theta_0$

Types of Regression Models

(Education)

x

Simple Linear regression model

y (Income)

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

weight for
the feature
feature

(Education)

x₁

Multiple Linear regression model

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$

(Soft Skills)

x₂

y (Income)

(Experience)

x₃

(Age)

x₄

Regression
Models

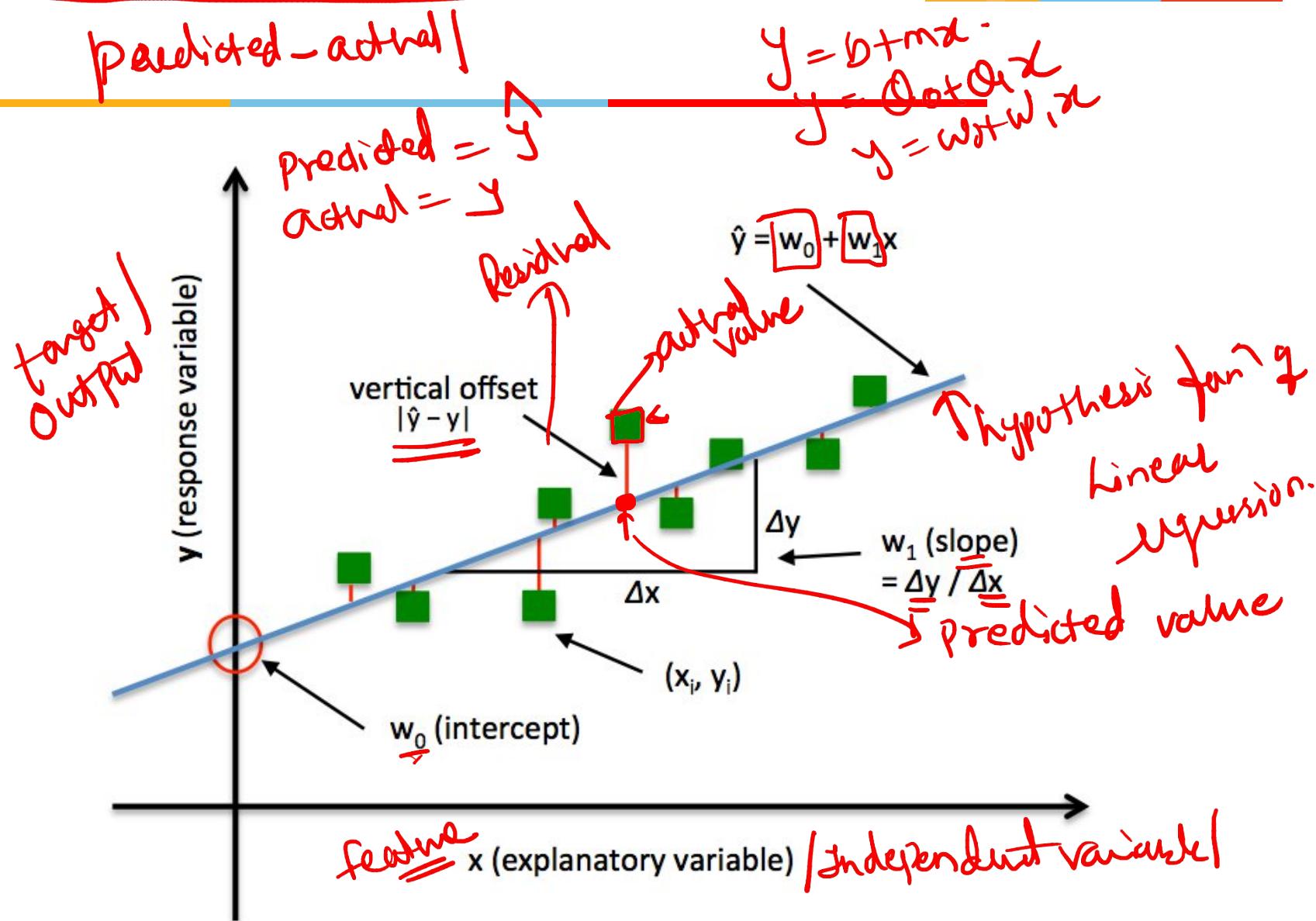


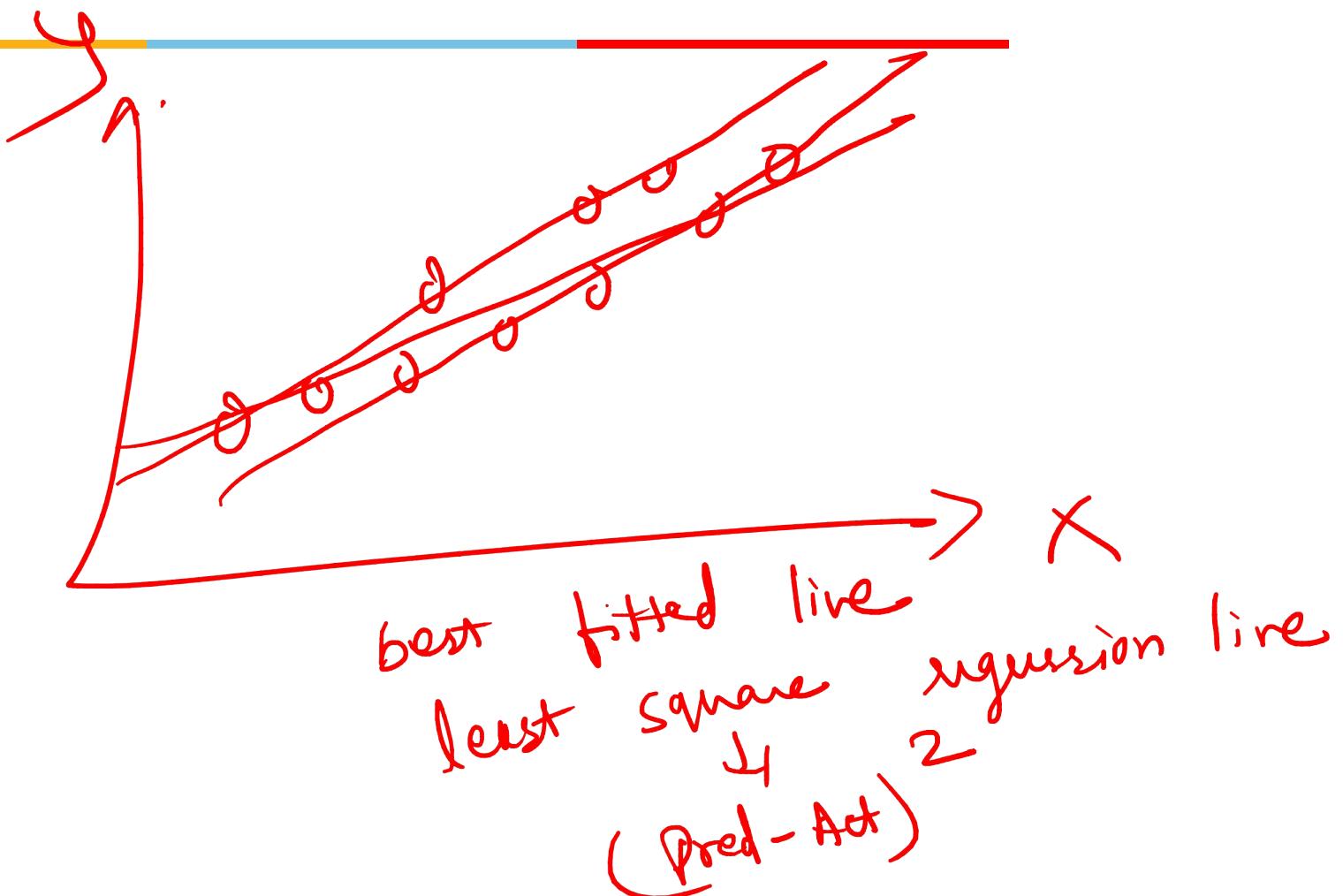
- Hypothesis:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d = \sum_{j=0}^d \theta_j x_j$$

Assume x₀ = 1

Least square regression line





Linear Regression – Hypothesis function

- Hypothesis:

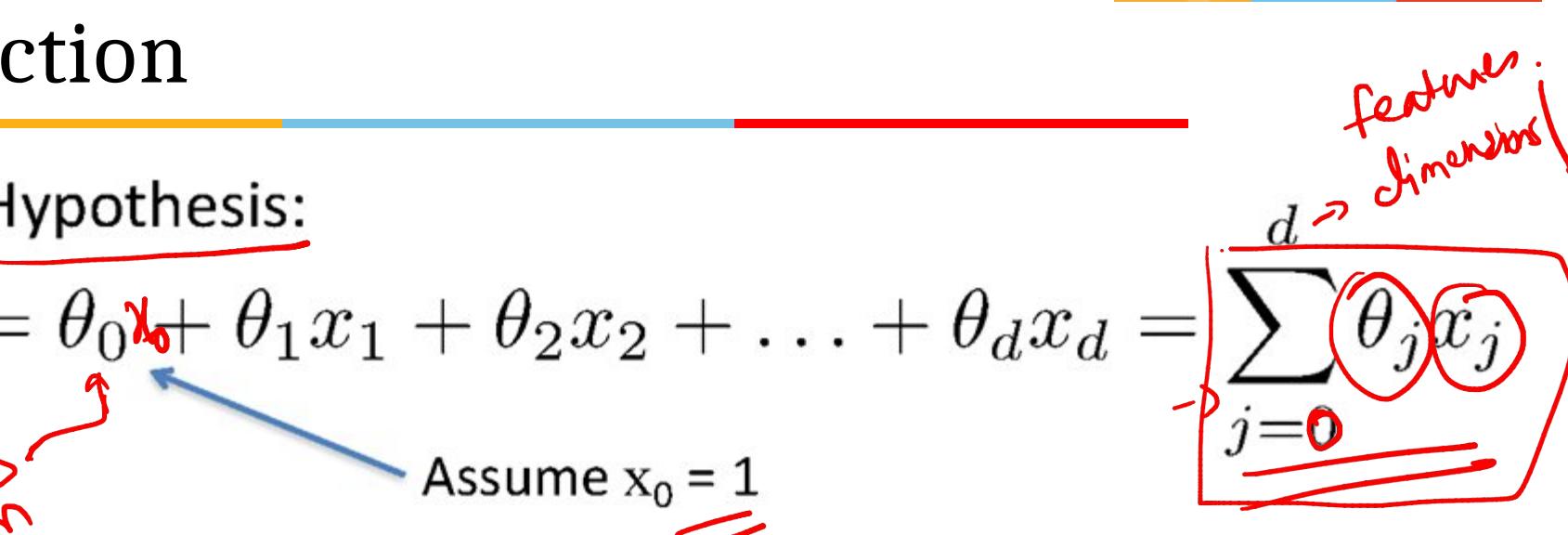
$$y = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d = \sum_{j=0}^d \theta_j x_j$$

d → features dimensions

θ₀ - θₖ → weights

bias term

Assume $x_0 = 1$



$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$

weights and bias ⇒ parameters

Learning the Model Parameters

Closed Form Solution Approach

Gradient Descent Approach

How to choose θ_i 's ?

Linear Regression : Least Squares

$$f(\vec{x})$$

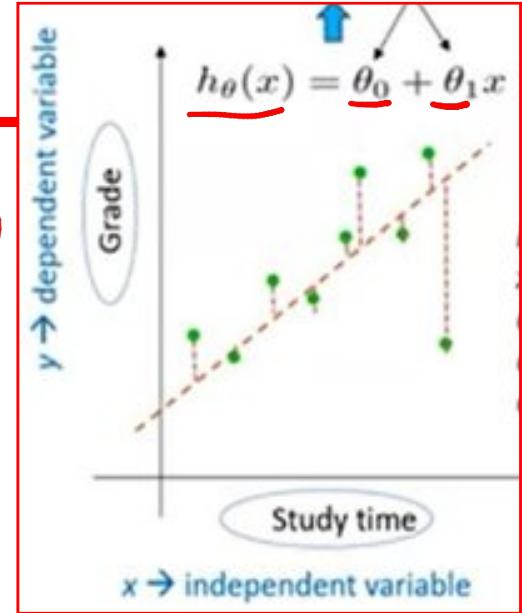
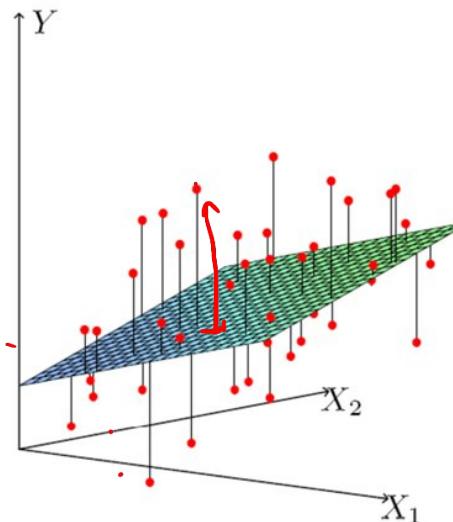
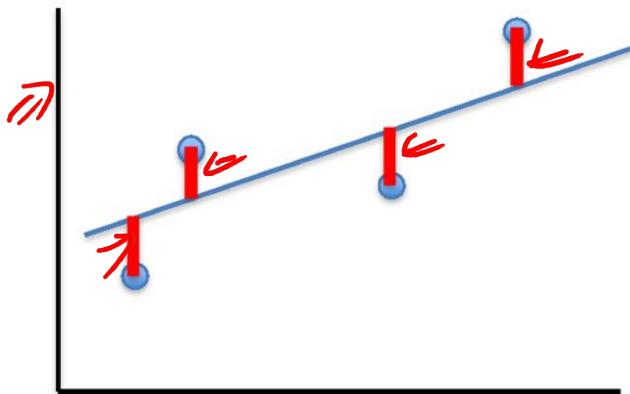
$n = \text{no. of training examples}$

- Cost Function

$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 \quad \begin{matrix} \text{pred} & - \text{Actu} \\ \text{(mean square error)} \end{matrix}$$

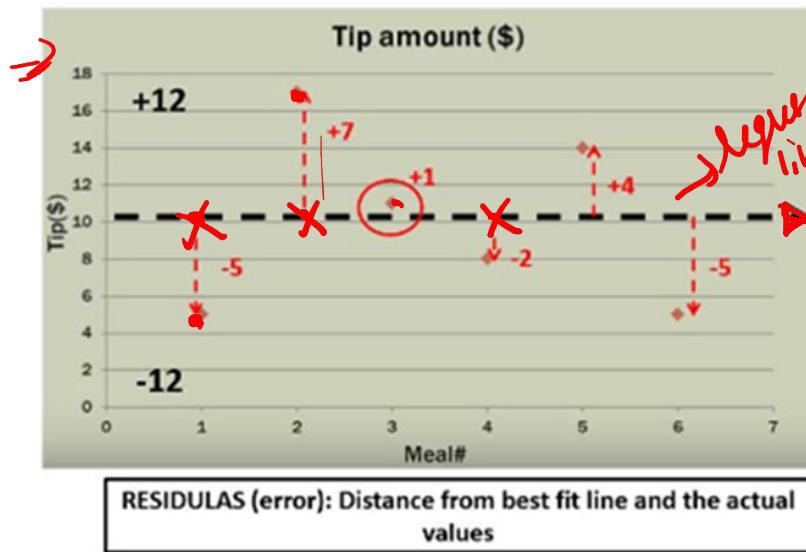
(MSE)

- Fit by solving $\min_{\theta} J(\theta)$



Notion of Cost Function

Meal #	Tip amount (\$)
1	5.00
2	17.00
3	11.00
4	8.00
5	14.00
6	5.00
mean = \$10	



pred - act

Meal#	Residual	Residual ²
1	-5	25
2	+7	49
3	+1	1
4	-2	4
5	+4	16
6	-5	25

Sum of Squared Errors (SSE) = 120

\hat{y}

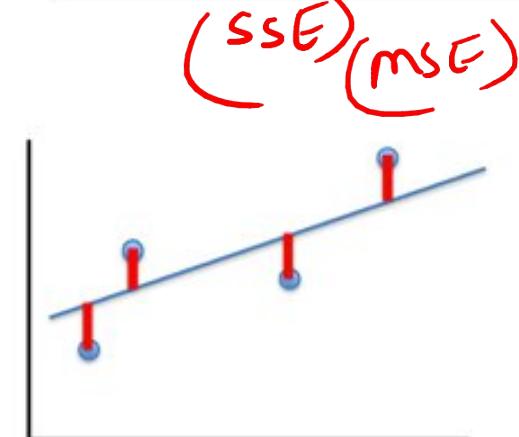
y

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

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

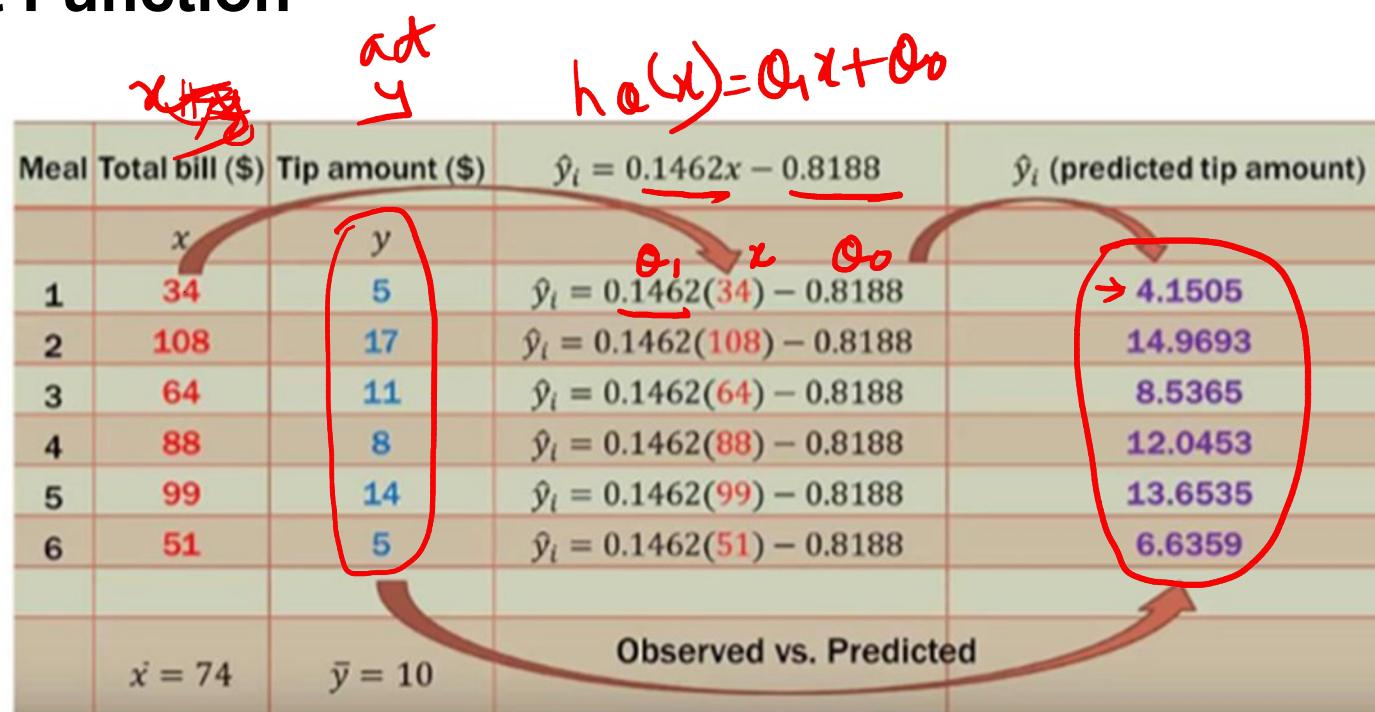
pred → \hat{y}

pred - act



Notion of Cost Function

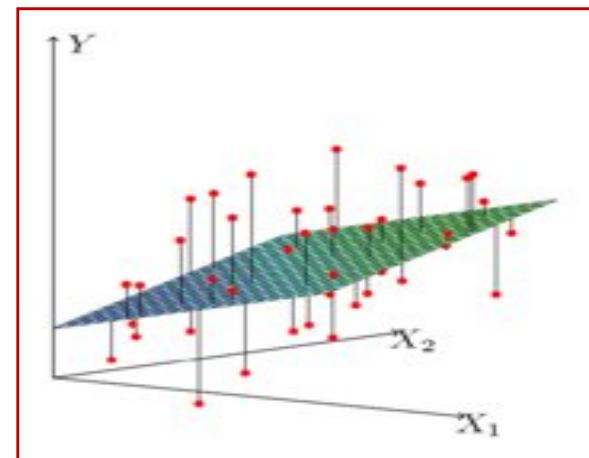
Meal #	Tip amount (\$)
1	5.00
2	17.00
3	11.00
4	8.00
5	14.00
6	5.00
mean = \$10	



$$y = \theta_0 + \theta_1 x$$

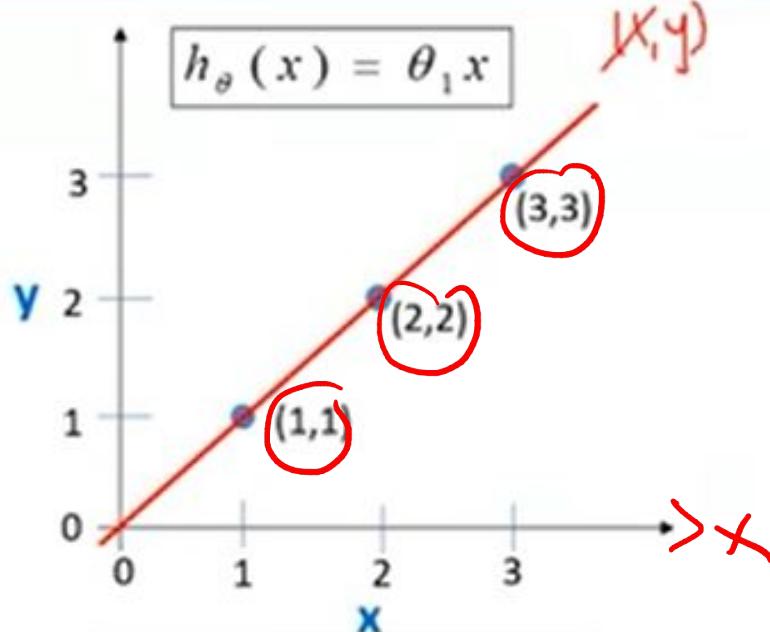
$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2$$

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



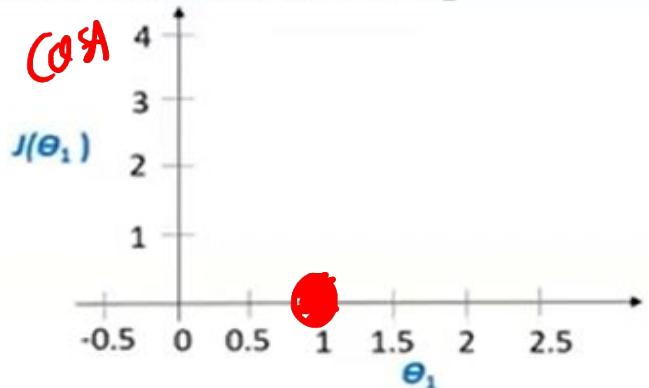
Understanding Cost Function

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



Cost Function

Function of parameter θ_1

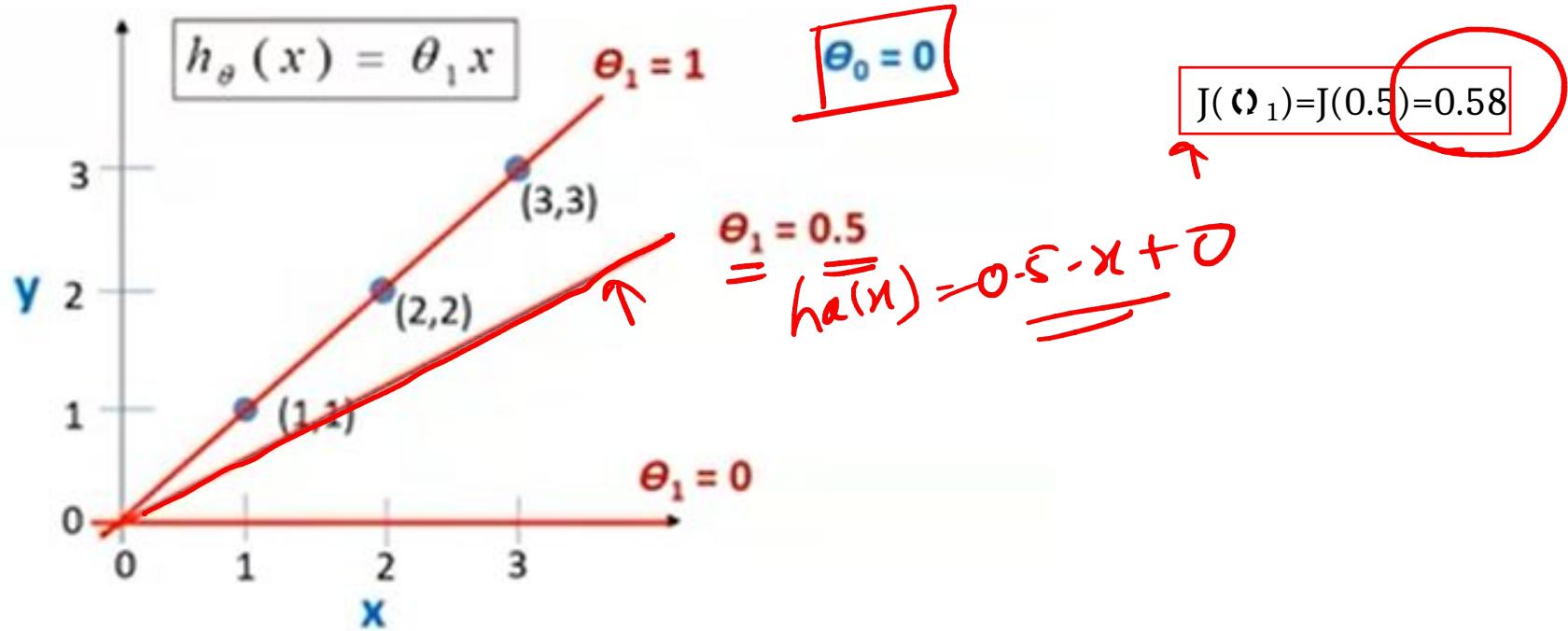


$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^n (\theta_1 x^{(i)} - y^{(i)})^2 \end{aligned}$$

$$\begin{aligned} J(\theta_1) &= \frac{1}{2 \cdot 3} ((1 \cdot 1 - 1)^2 + (1 \cdot 2 - 2)^2 + (1 \cdot 3 - 3)^2) \\ &= 0 \end{aligned}$$

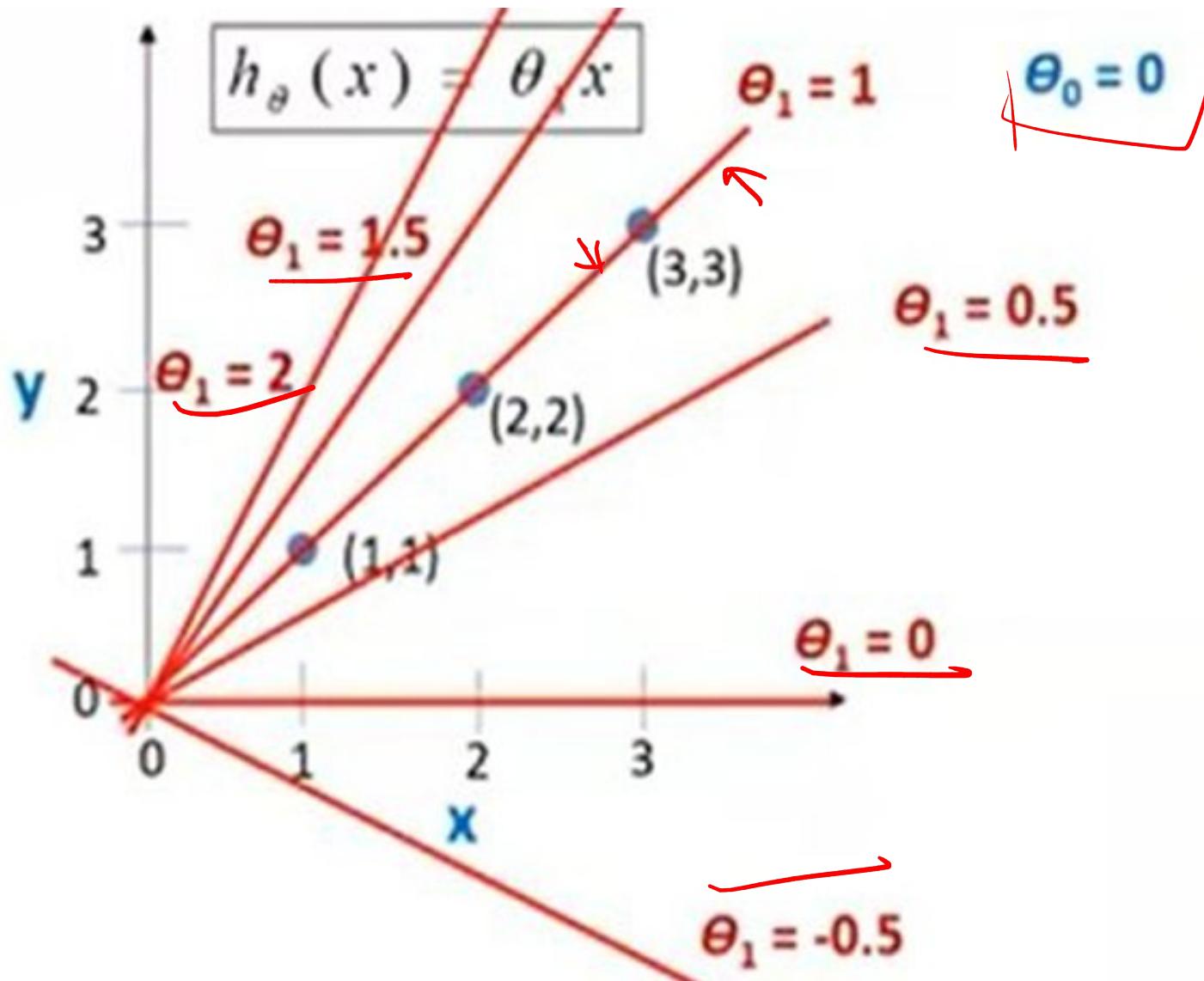
$n = 3$ = no. of training examples

Understanding Cost Function



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

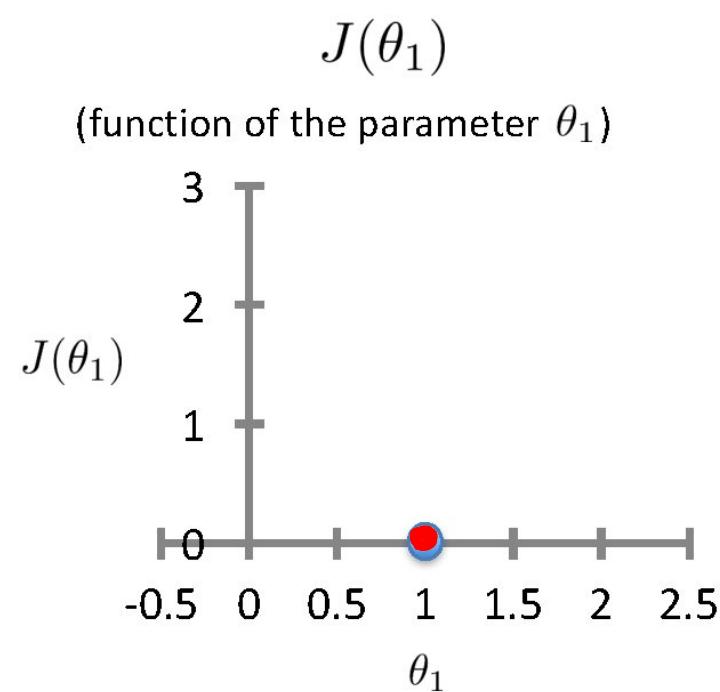
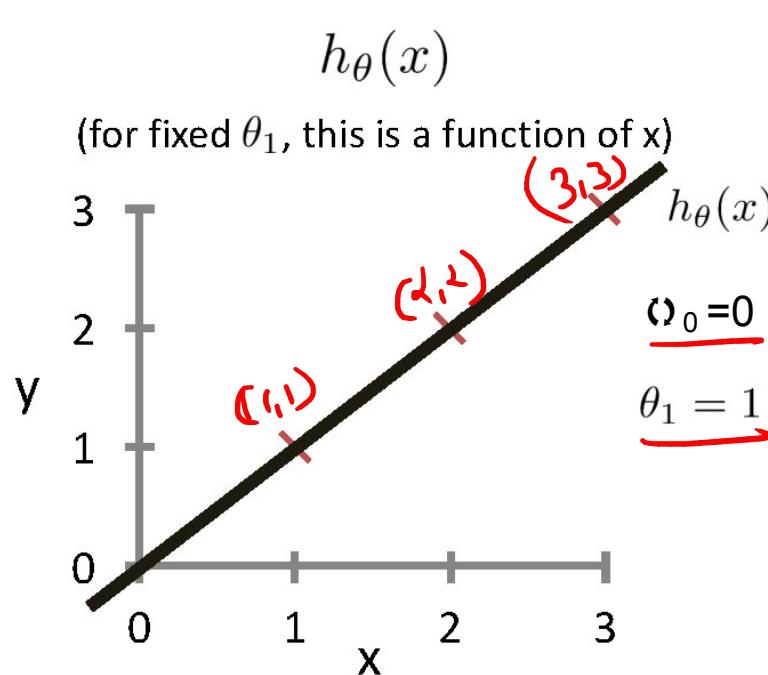
Understanding Cost Function



Intuition Behind Cost Function

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

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$



Based on example
by Andrew Ng



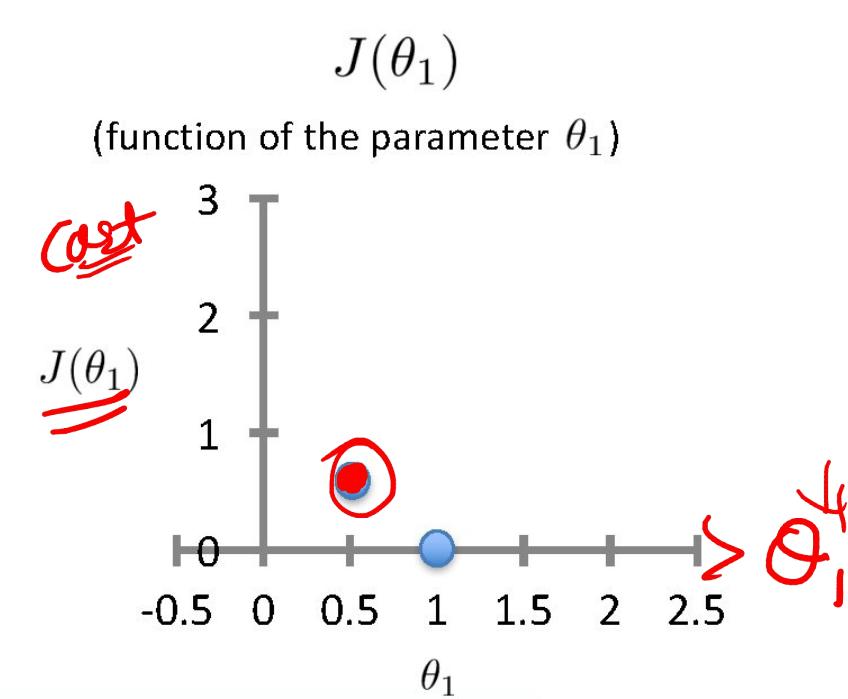
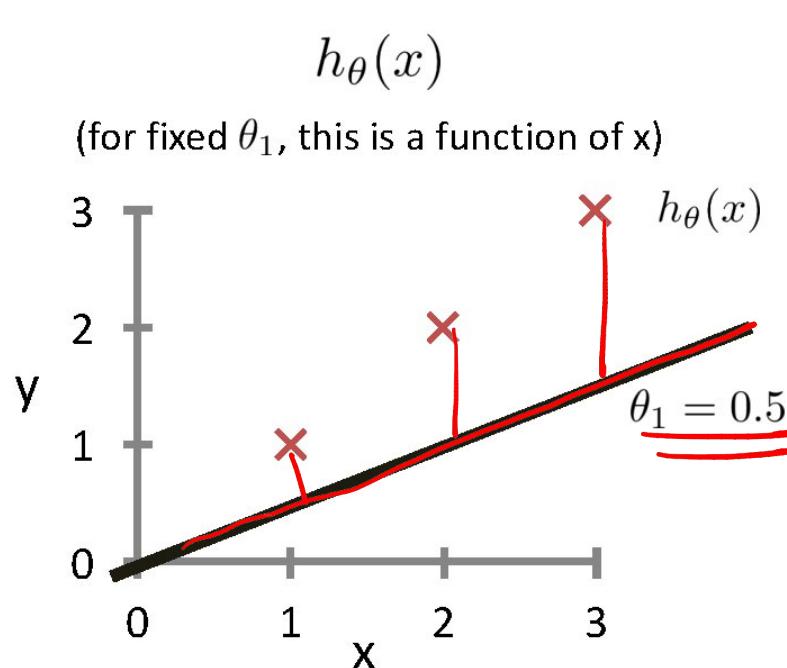
lead

Intuition Behind Cost Function

$$h_{\theta}(x) = \sum \theta_i x^i$$

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

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$



Based on example
by Andrew Ng

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



7

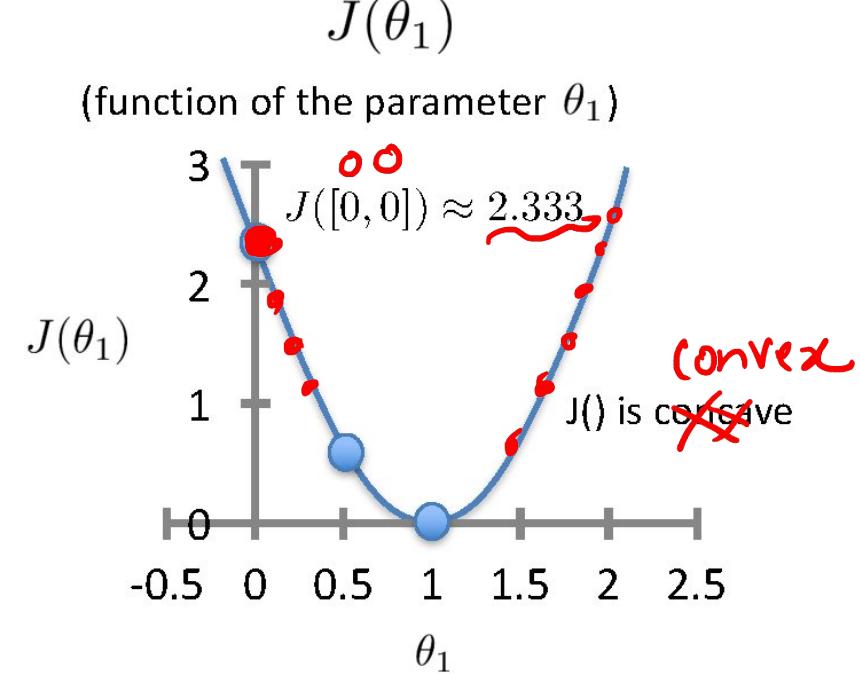
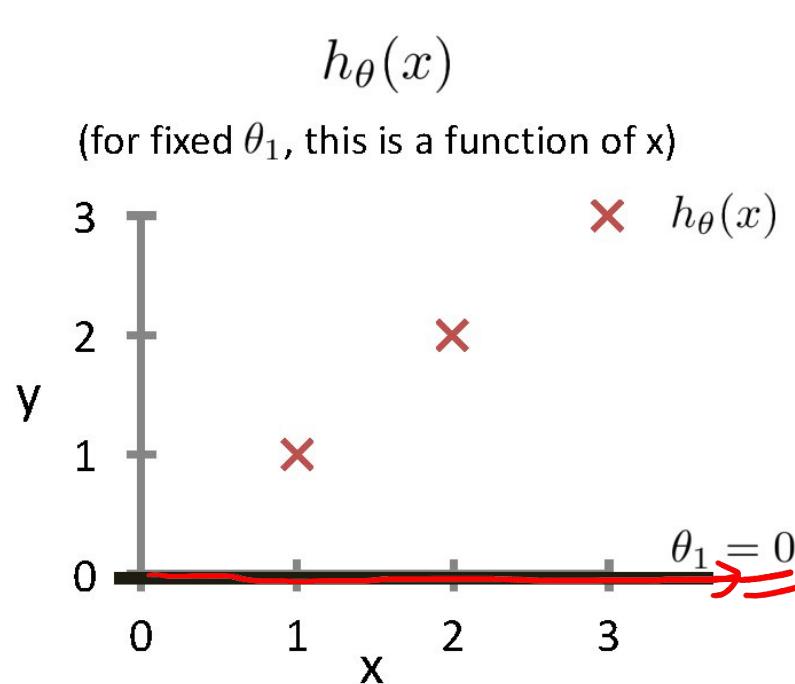
$\theta_0 \Rightarrow 1$
 $\theta_1 \Rightarrow 0.5$
 $\theta_2 \Rightarrow 0.4$

Intuition Behind Cost Function

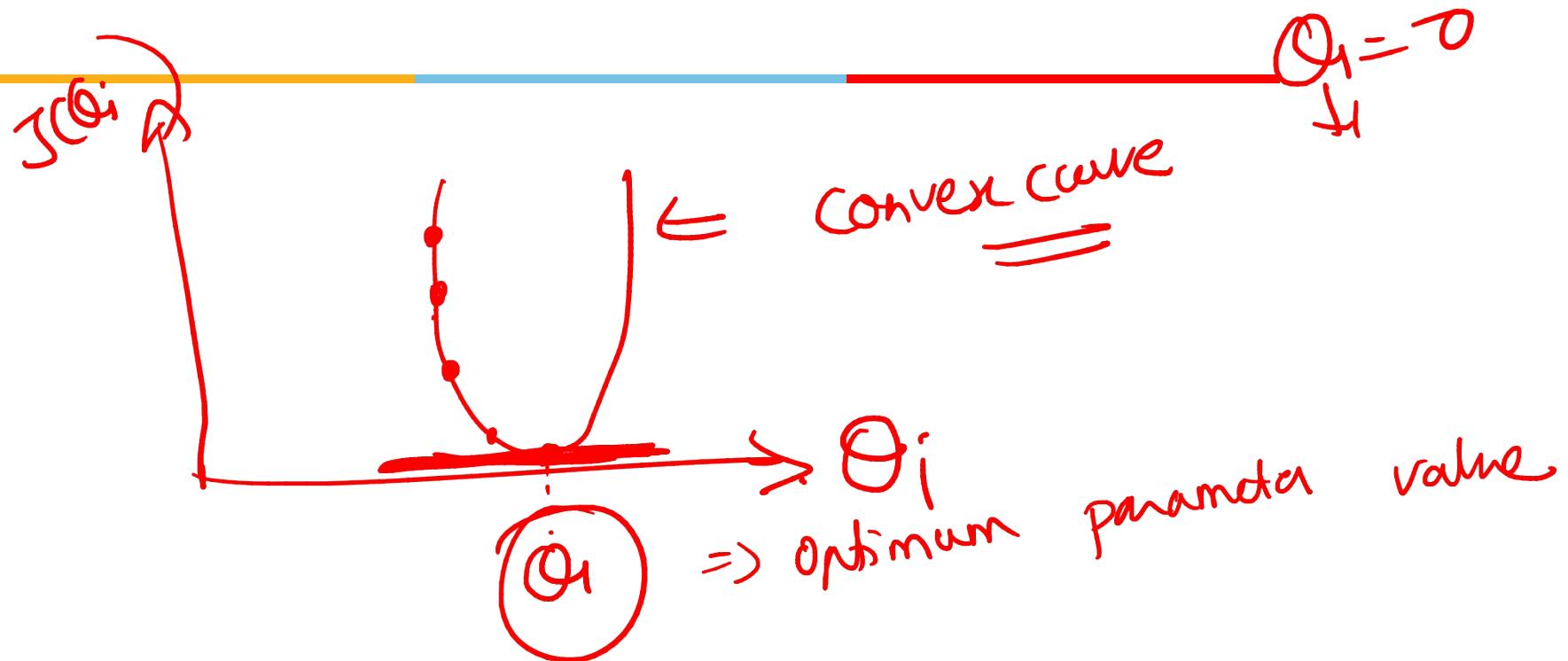
$$J(\theta) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

 \rightarrow concave
 \rightarrow convex

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$



Based on example
by Andrew Ng



Linear Regression Summary



weights over features.

Hypothesis function

$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$

Parameters to learn

$$\theta_i$$

Cost function

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

act - pred. ms^E

Convex

Concave

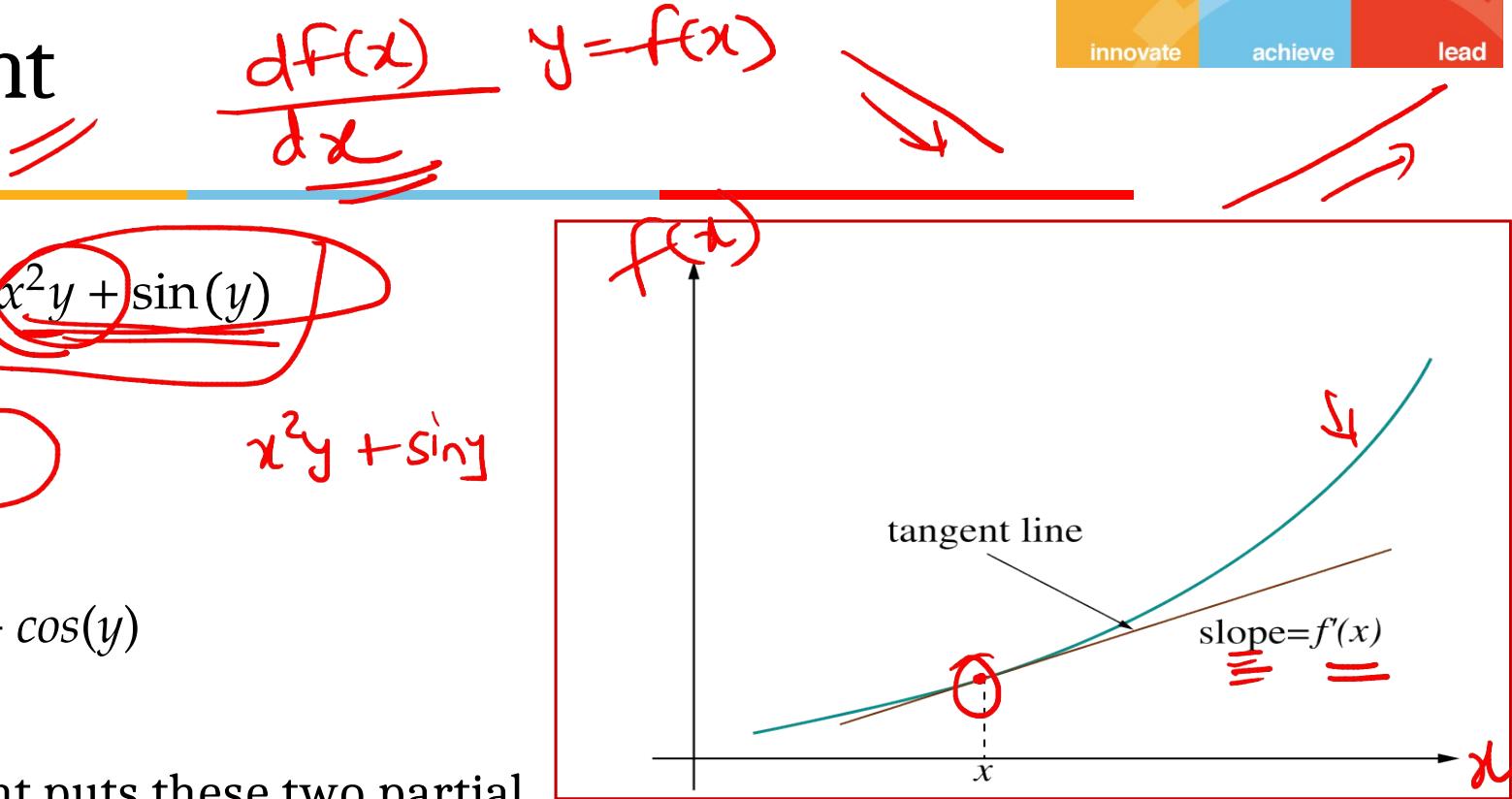
Goal

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

y

Gradient Descent

Gradient



$$\nabla f(x, y) = \nabla(x^2y + \sin(y)) = \begin{bmatrix} 2xy \\ x^2 + \cos(y) \end{bmatrix}$$

$$\frac{dy}{dx}$$

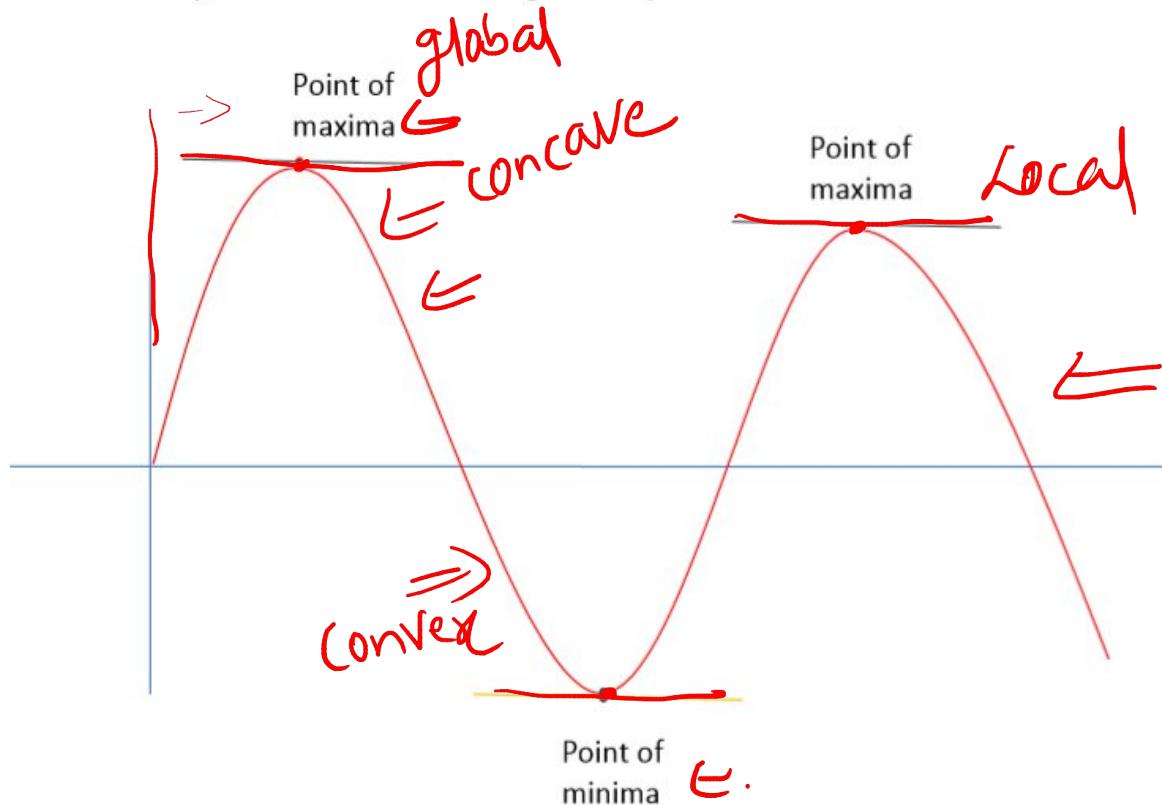
$$\nabla f(x, y) = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix}$$

Maxima and Minima

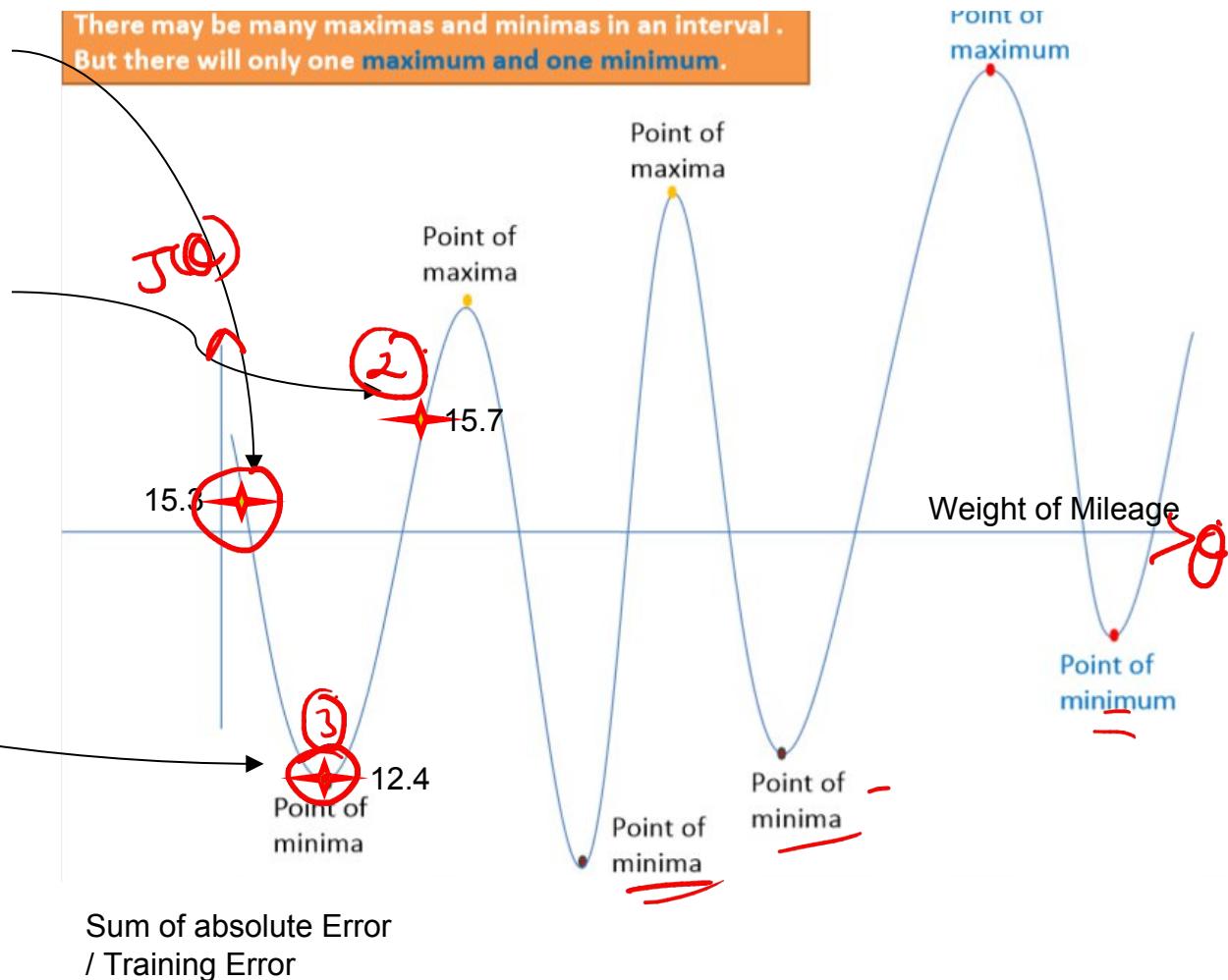
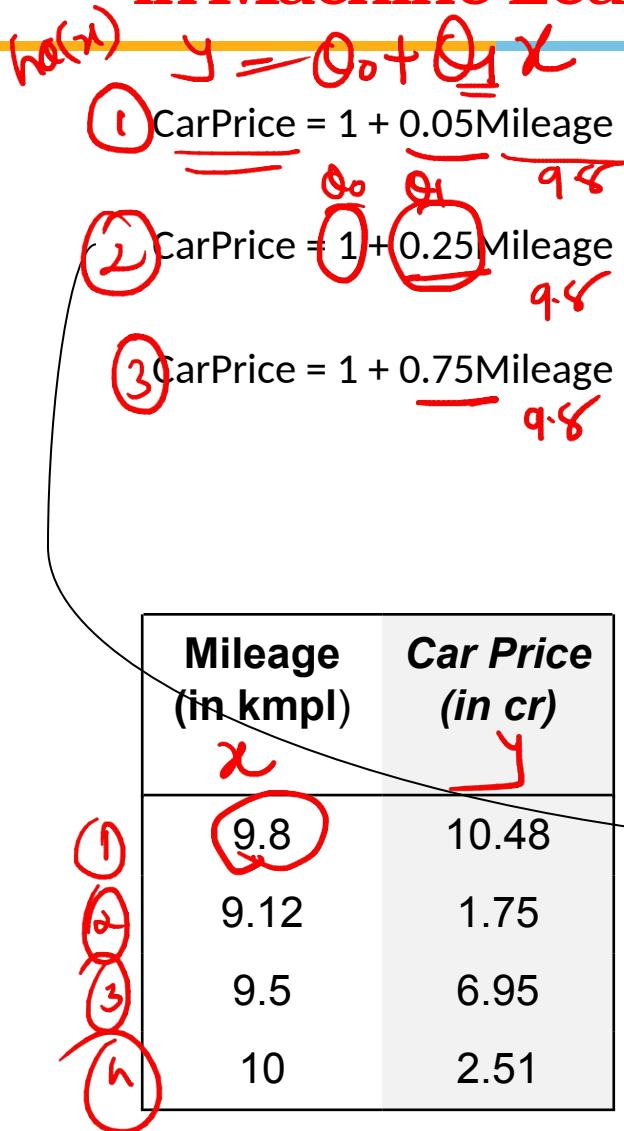
- For maxima and minima $m = dy/dx = \tan 0^\circ = 0$
- $dy/dx = 0$ means tangent is parallel to X-axis.

local minima
global minima

$$f'(x) = 0$$



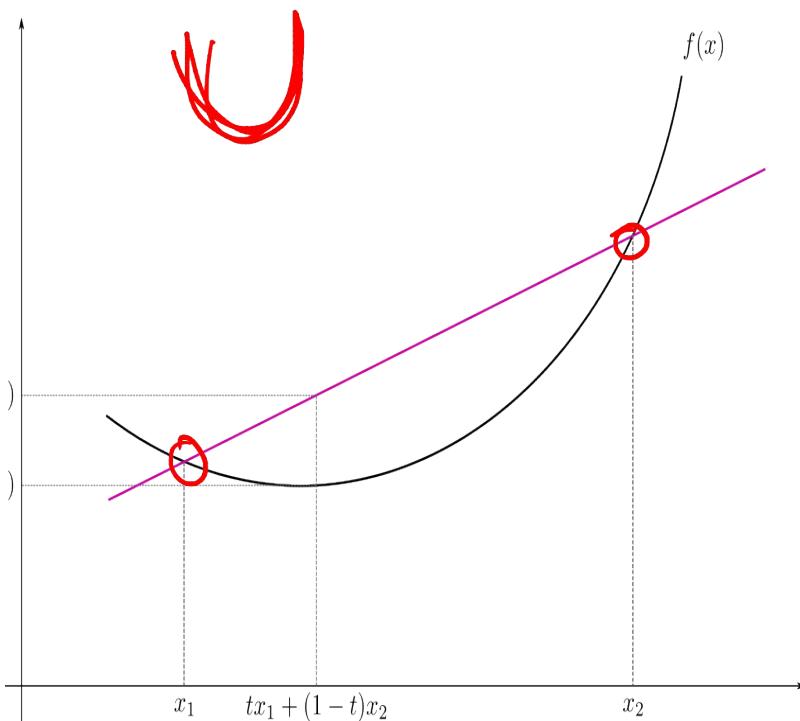
Notion of Maxima and Minima of a function in Machine Learning



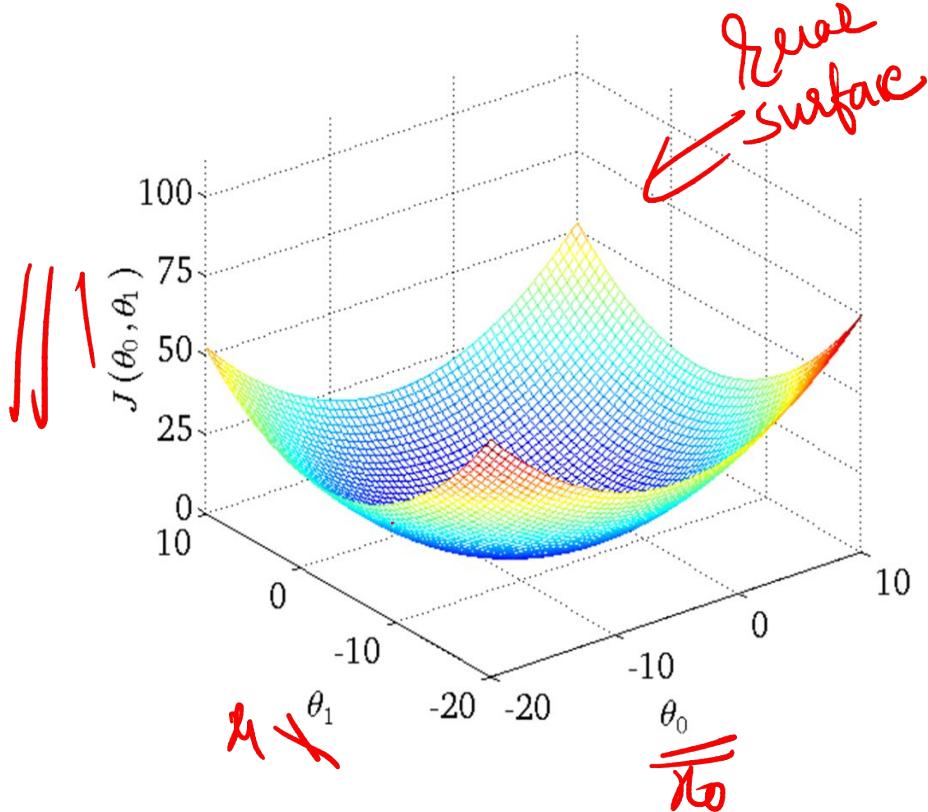
Convex Function



Univariate



Multivariate



Real-valued function defined on an n -dimensional interval is called **convex** if the line segment between any two points on the graph of the function lies above or on the graph

Basic Search Procedure

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

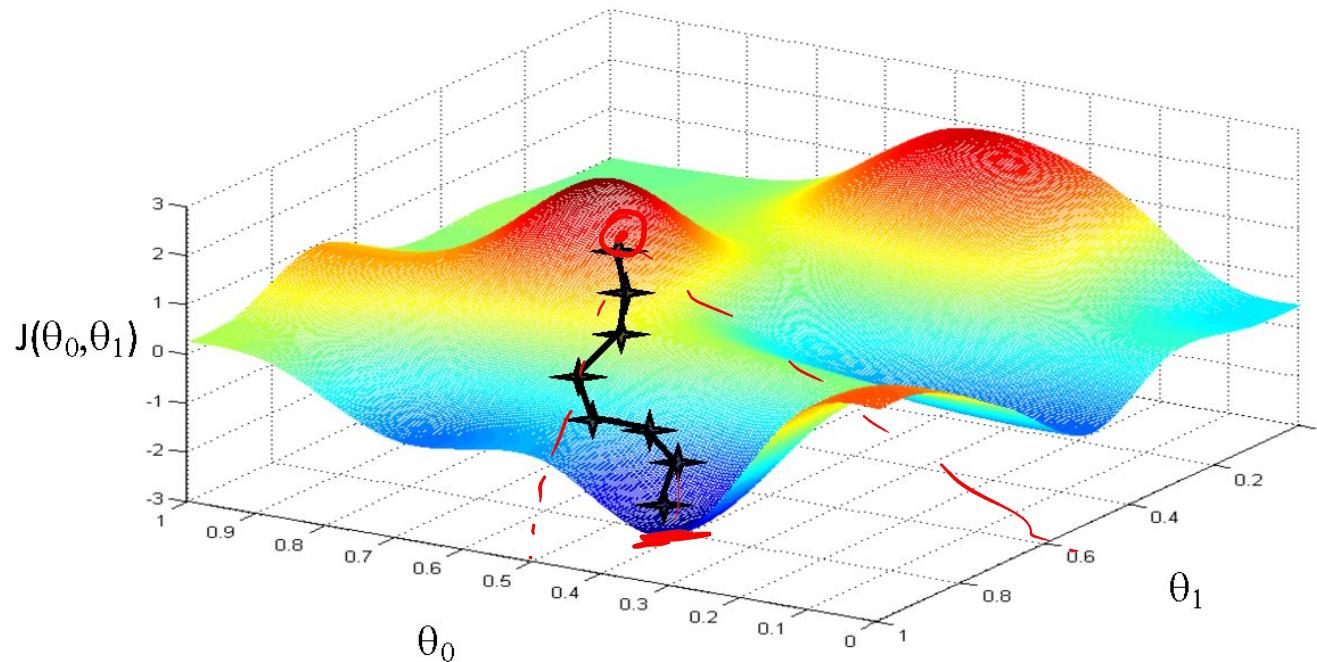


Figure by Andrew Ng

Basic Search Procedure

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

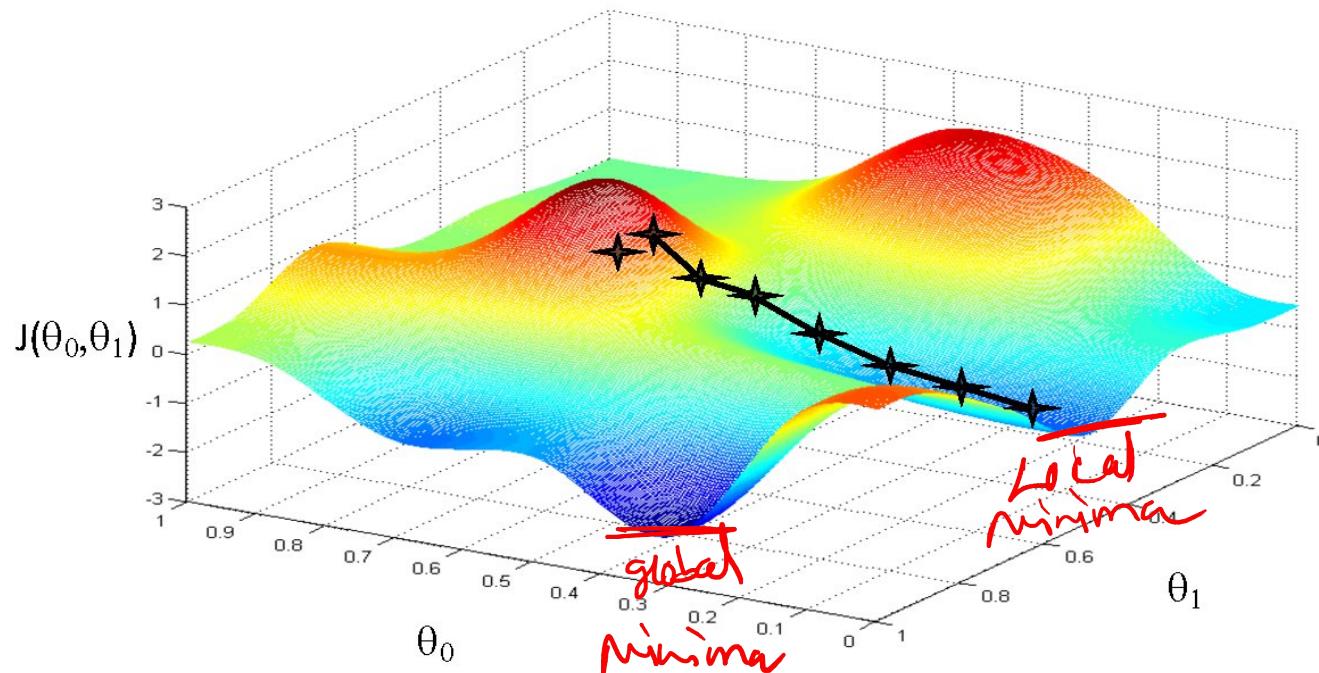
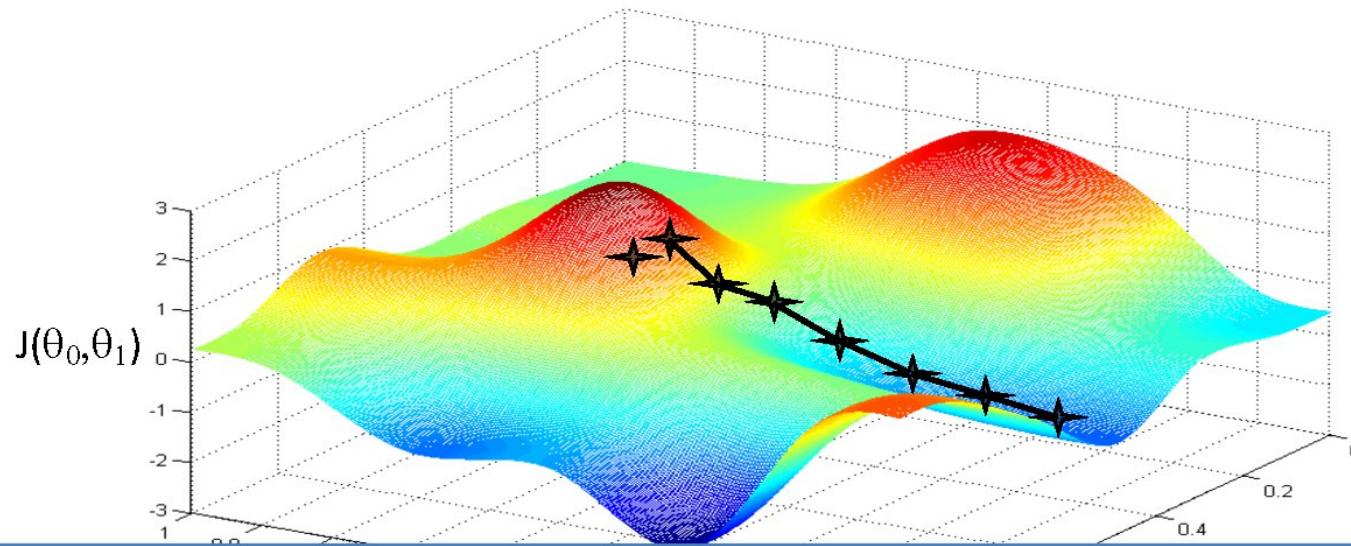


Figure by Andrew Ng

Basic Search Procedure

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



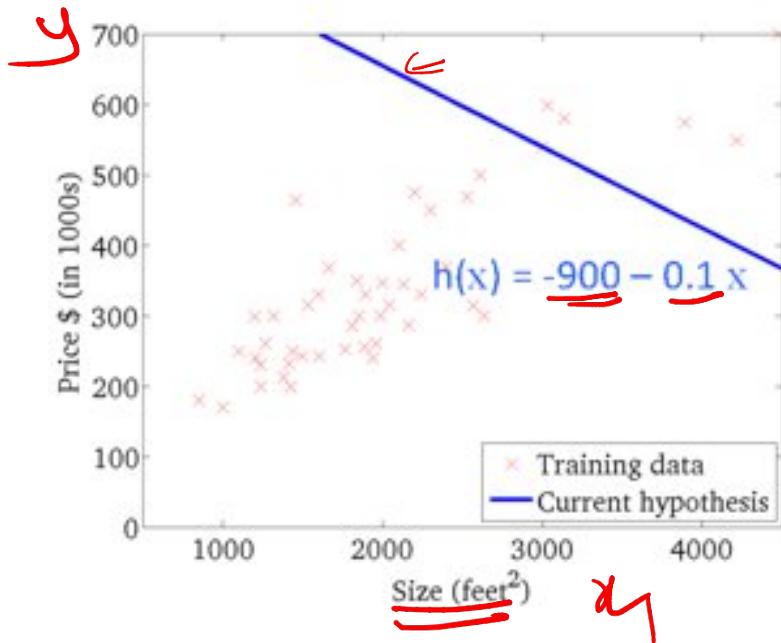
Since the least squares objective function is convex (concave),
we don't need to worry about local minima

Gradient Descent

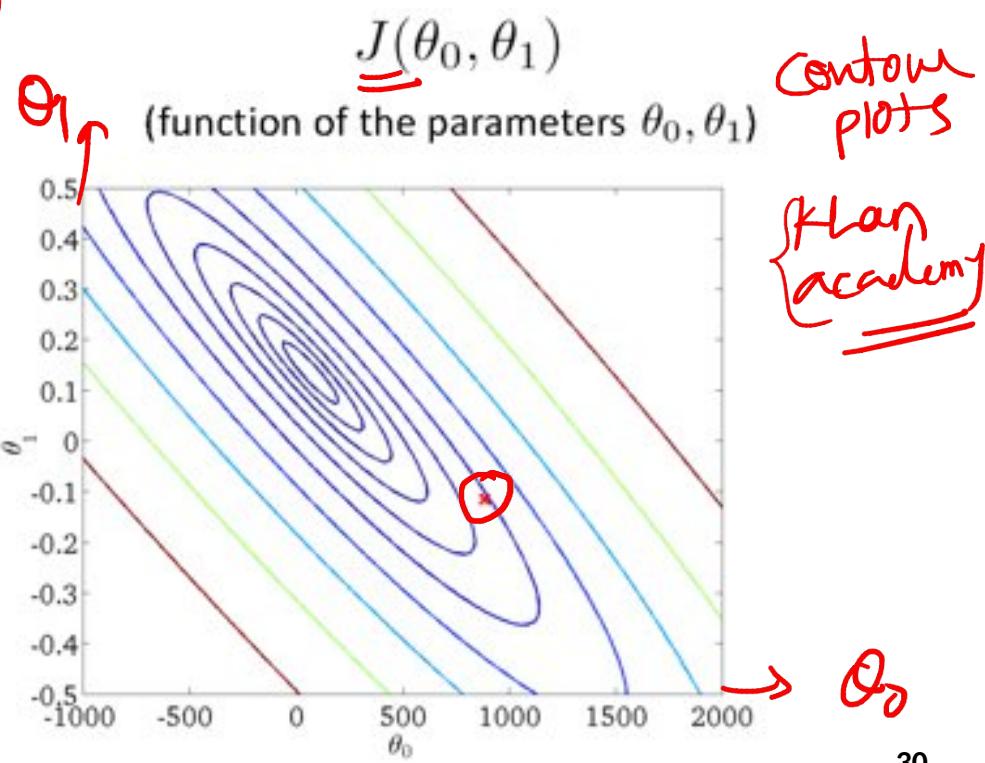
$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$

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

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

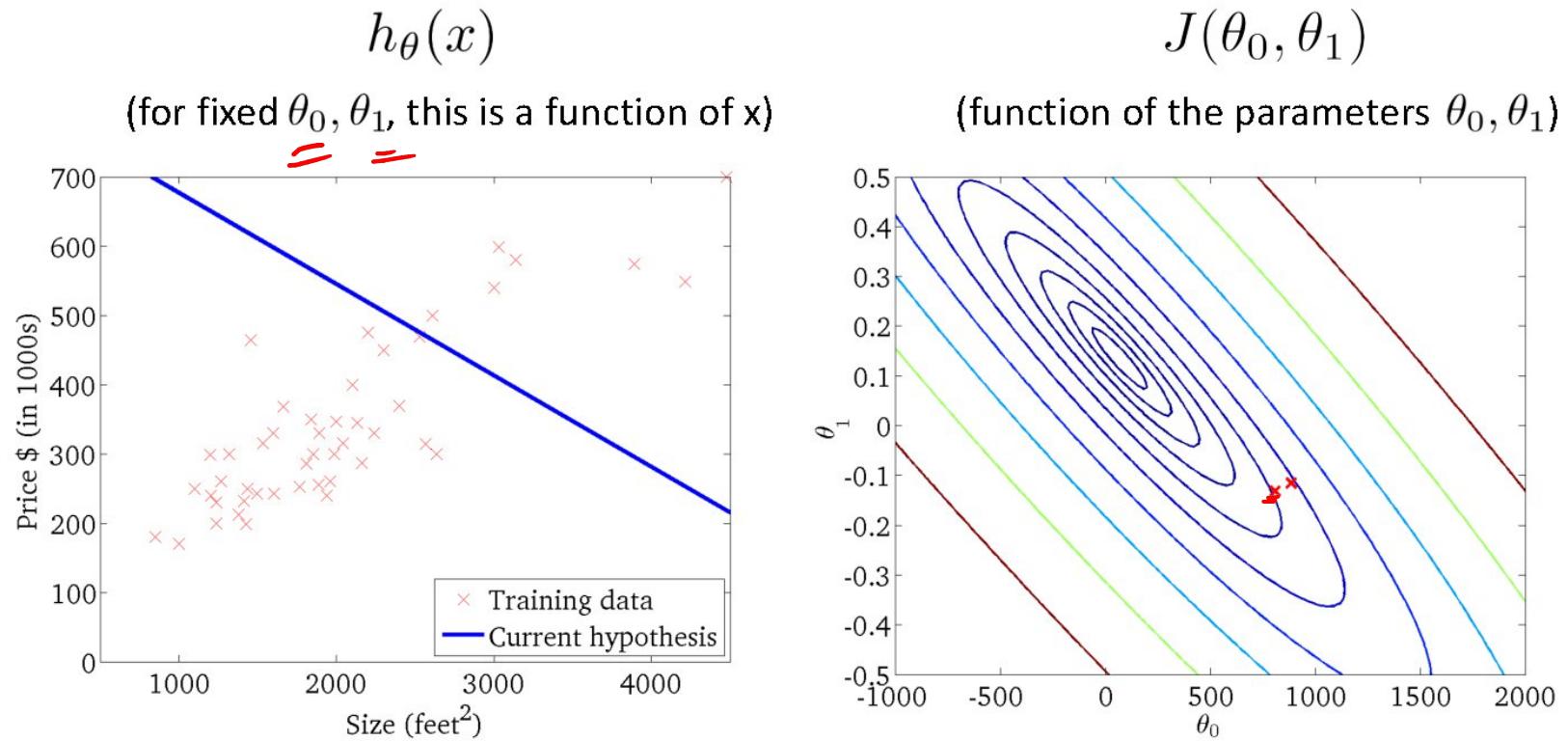


Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...





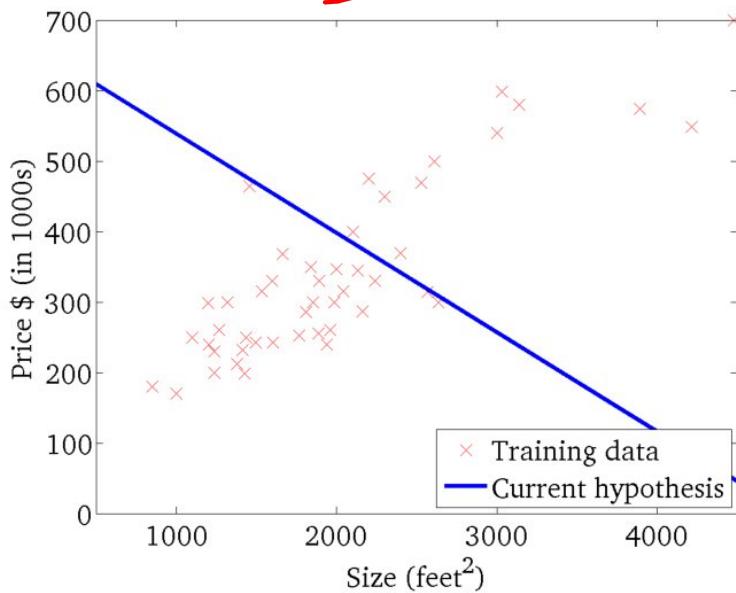
Gradient Descent



Gradient Descent

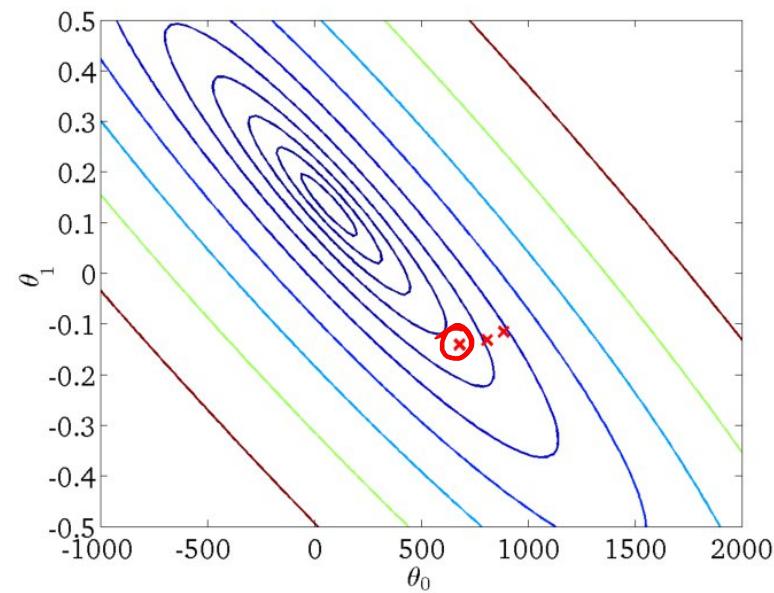
$$h_{\theta}(x)$$

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



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

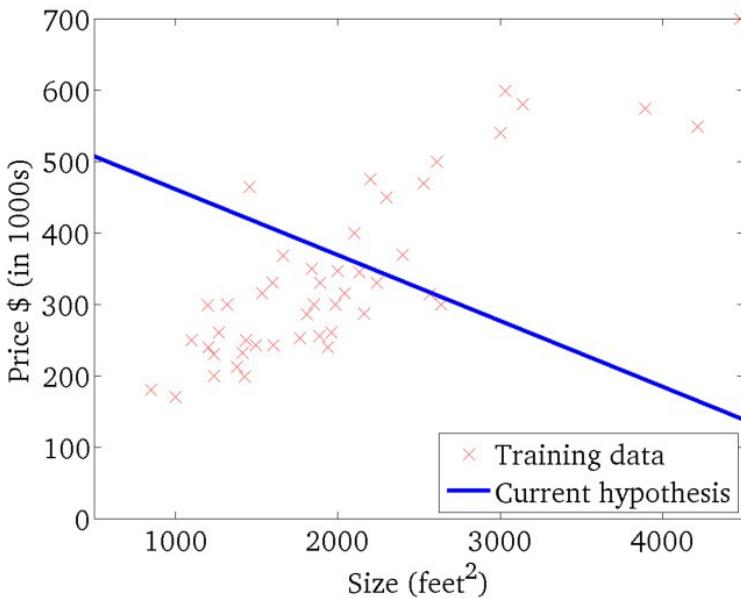
(function of the parameters θ_0, θ_1)



Gradient Descent

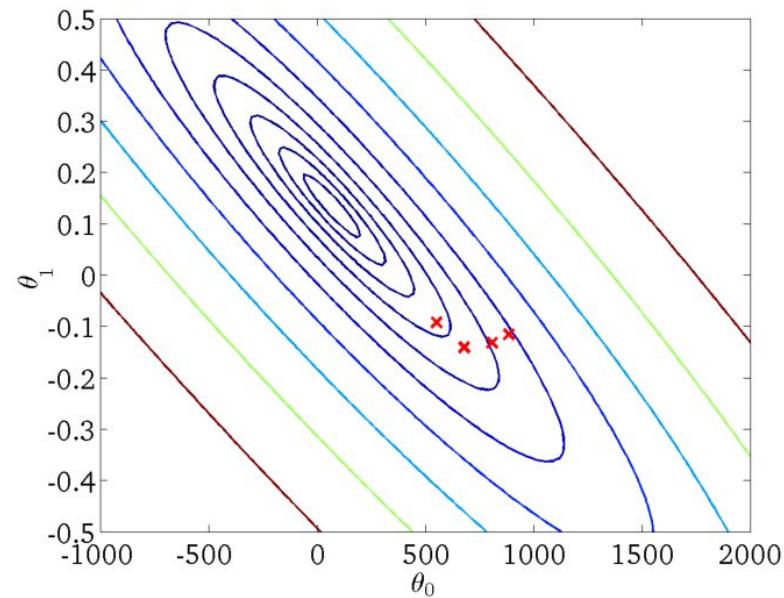
$$h_{\theta}(x)$$

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

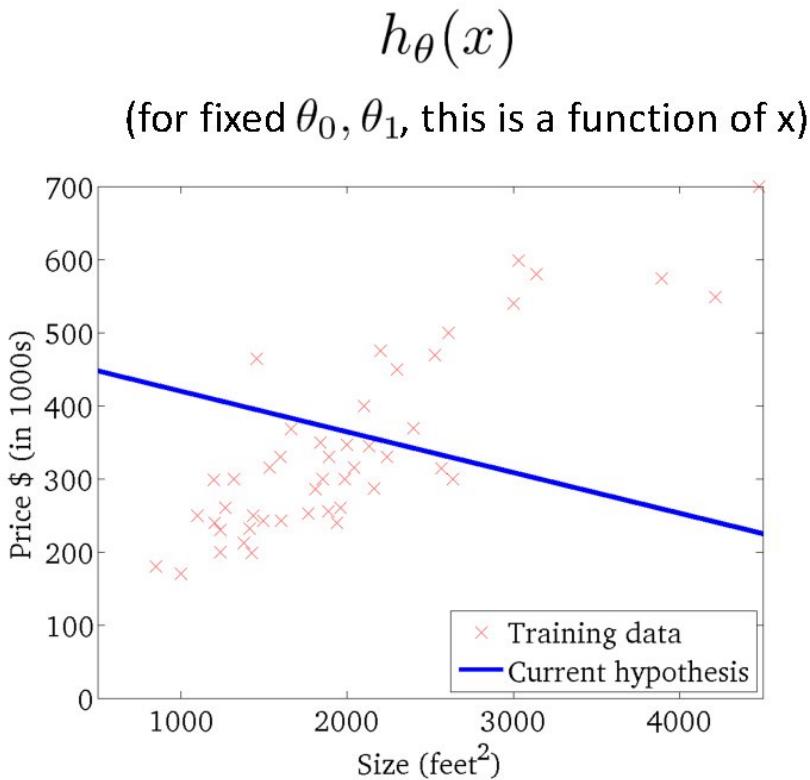


$$J(\theta_0, \theta_1)$$

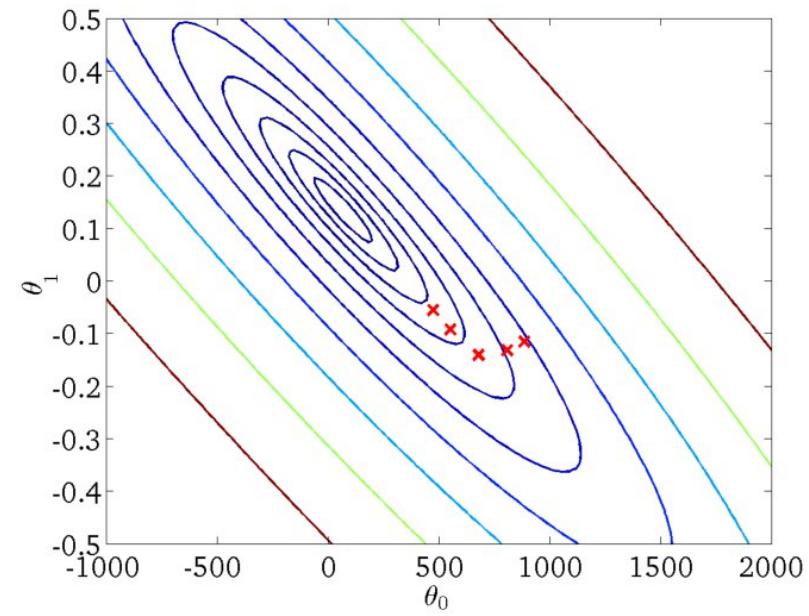
(function of the parameters θ_0, θ_1)



Gradient Descent



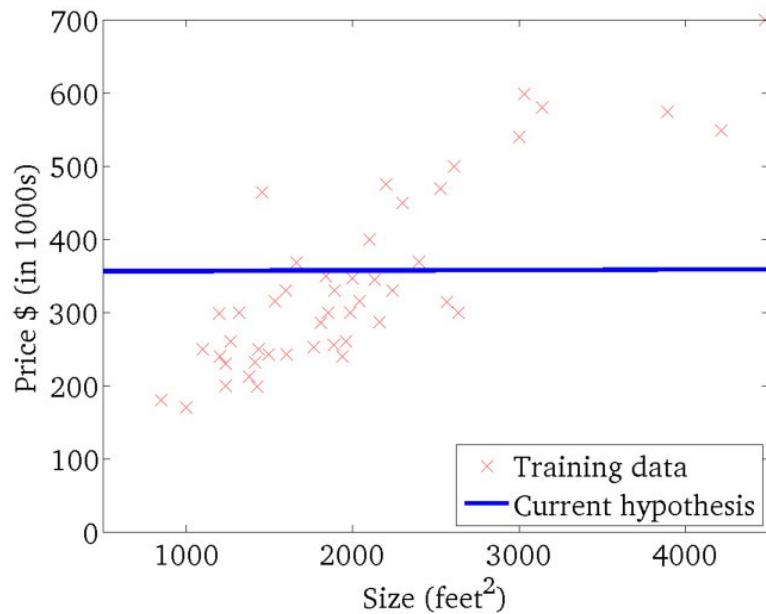
$J(\theta_0, \theta_1)$
(function of the parameters θ_0, θ_1)



Gradient Descent

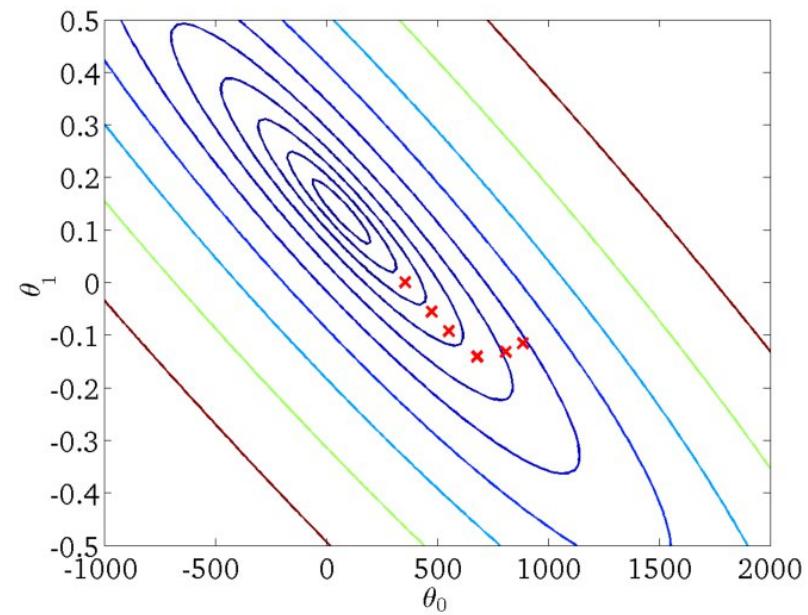
$$h_{\theta}(x)$$

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



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

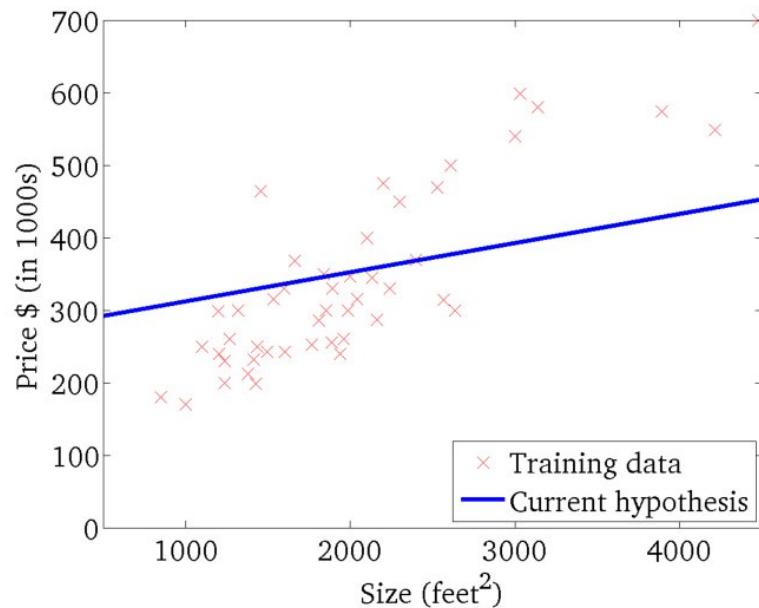
(function of the parameters θ_0, θ_1)



Gradient Descent

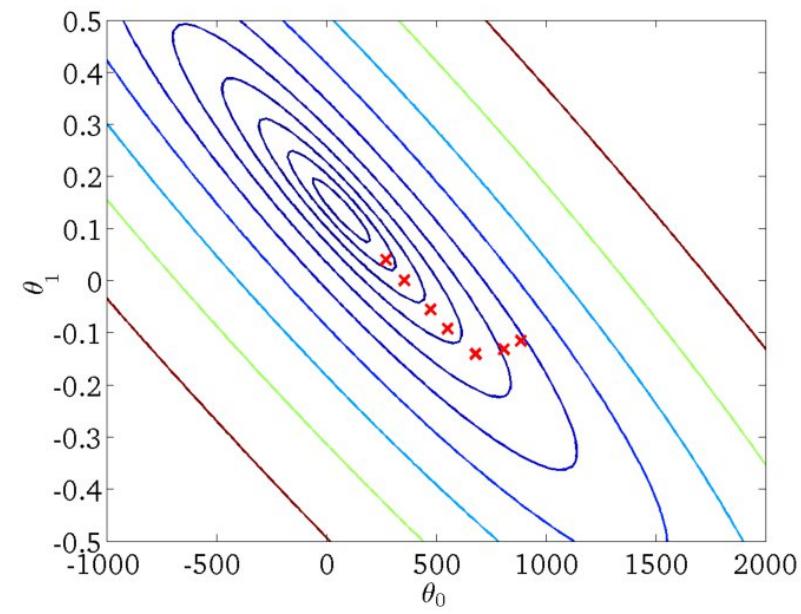
$$h_{\theta}(x)$$

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



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

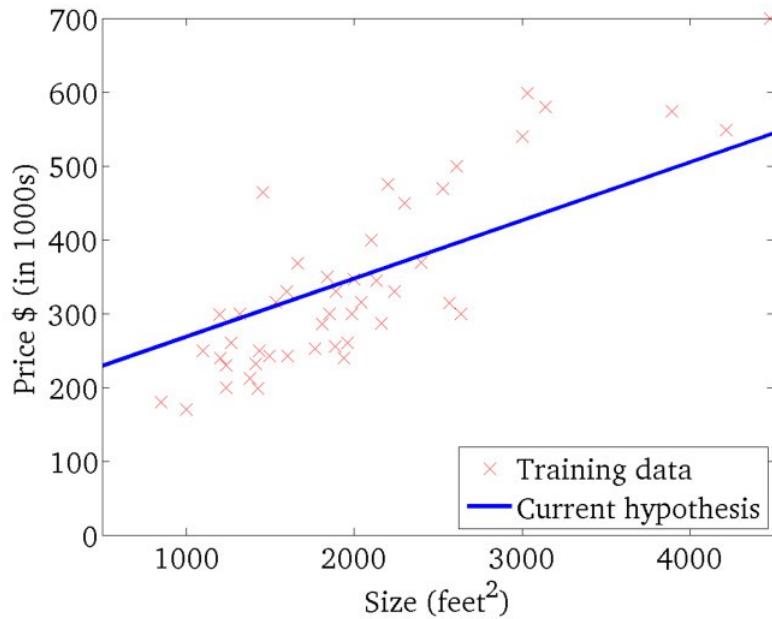
(function of the parameters θ_0, θ_1)



Gradient Descent

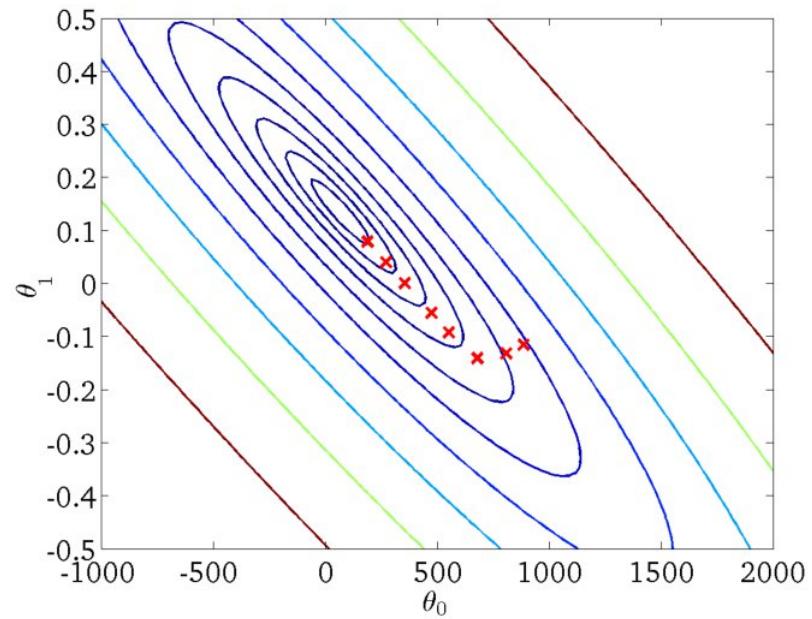
$$h_{\theta}(x)$$

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



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

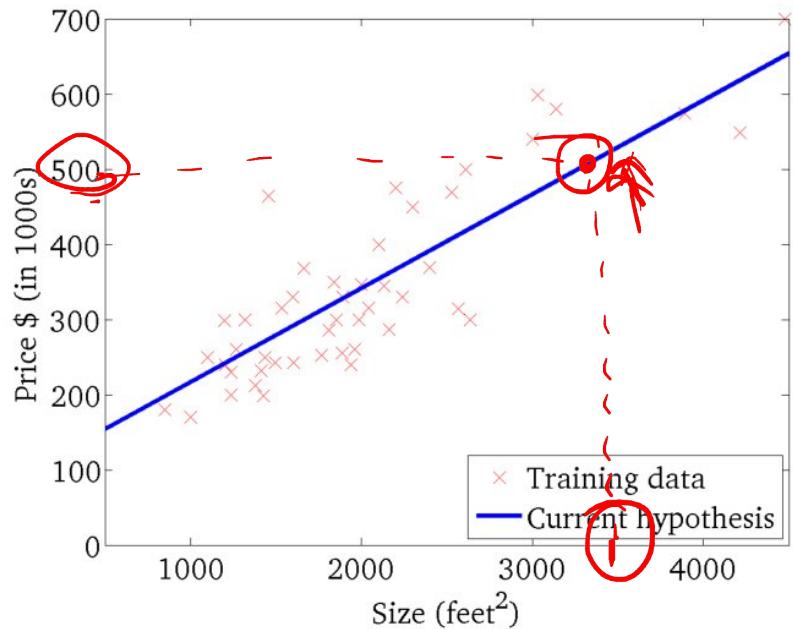
(function of the parameters θ_0, θ_1)



Gradient Descent

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

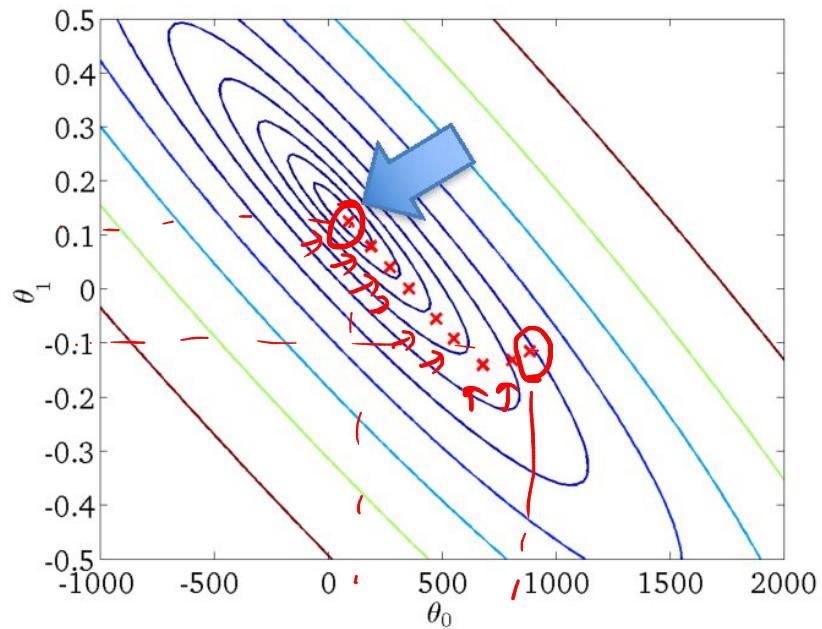
(for fixed θ_0, θ_1 , this is a function of x)



x

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

(function of the parameters θ_0, θ_1)



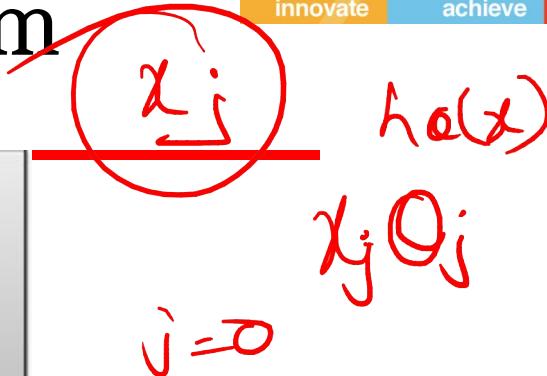
Gradient Descent Algorithm

- Initialize θ
- Repeat until convergence

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

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

simultaneous update
for $j = 0 \dots d$

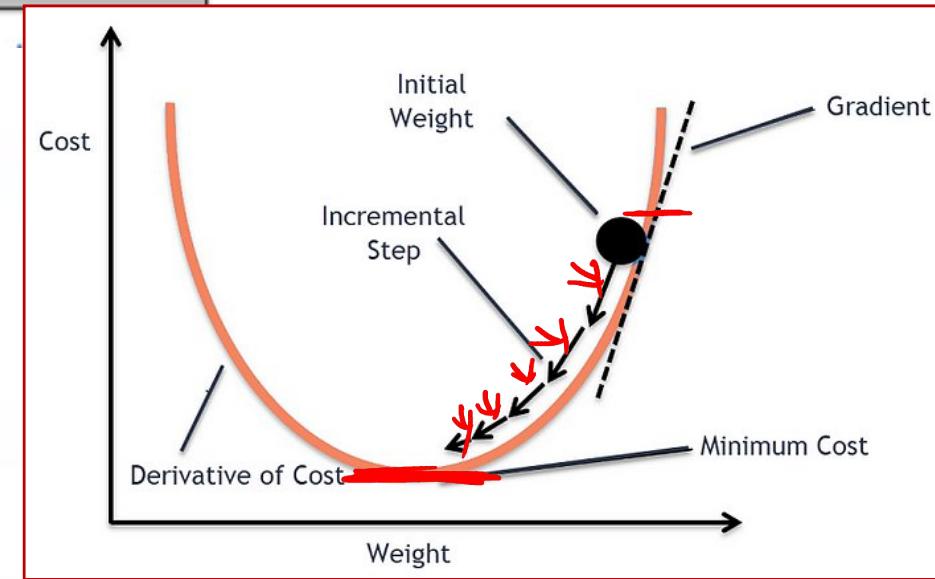


Have some function $J(\theta_0, \theta_1)$

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

Outline:

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



$$\sum_{j=0}^d x_j Q_j \quad \Rightarrow \quad j = 0 \dots d$$

$$Q_0^{new} = Q_0^{old} - \alpha \frac{\partial J(Q)}{\partial Q_0}$$

$$Q_1^{new} = Q_1^{old} - \alpha \frac{\partial J(Q)}{\partial Q_1}$$

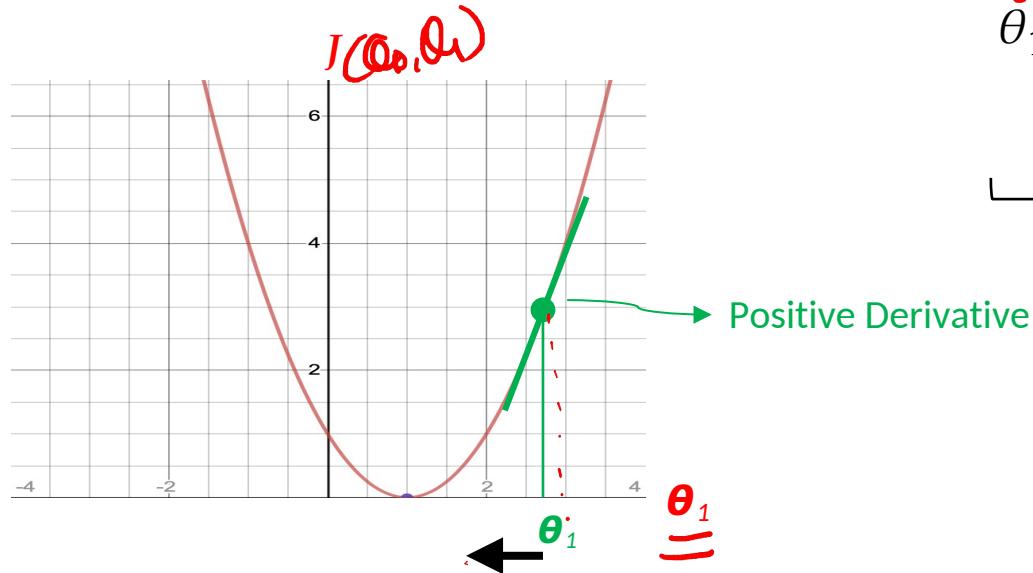
$$\vdots$$

$$Q_d^{new} = Q_d^{old} - \alpha \frac{\partial J(Q)}{\partial Q_d}$$

$$Q_j^{new} = Q_j^{old} - \alpha \frac{\partial J(Q)}{\partial Q_j}$$

Gradient Descent Intuition

- optimization objective is to minimize $J(\theta_1)$



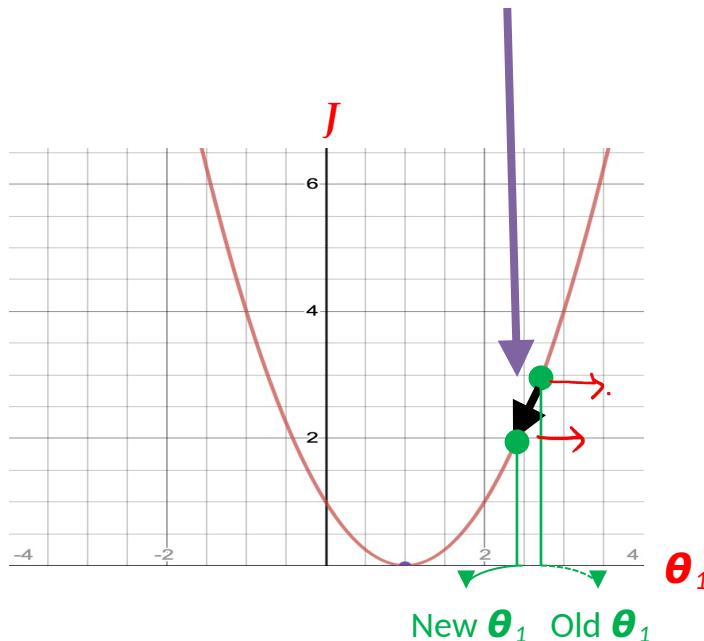
$$\begin{aligned} \text{new } \theta_1 &= \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j} \\ &= \theta_1 - \alpha \text{ (Positive Number)} \end{aligned}$$

Decrease $\underline{\theta_1}$ by a certain value

The Impact of Partial Derivative

- optimization objective is to minimize $J(\theta_1)$
 θ_0, θ_1

Step towards negative direction of a gradient

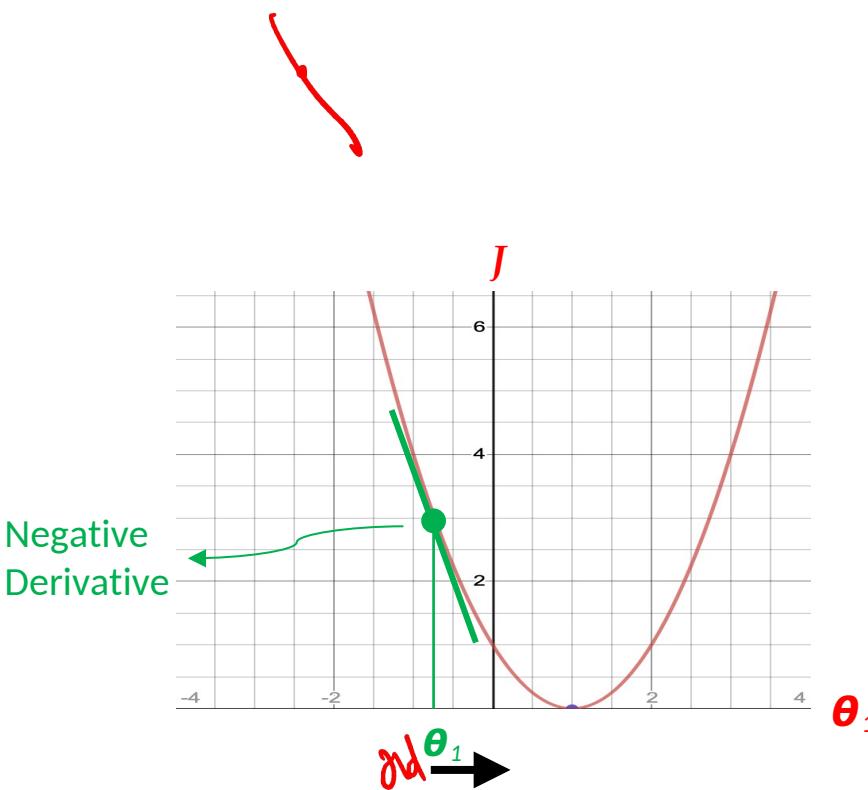


$$\begin{aligned}
 \theta_1 &= \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j} \\
 &= \theta_1 - \alpha (\text{Positive Number})
 \end{aligned}$$

Decrease θ_1 by a certain value

The Impact of Partial Derivative

- optimization objective is to minimize $J(\theta_1)$
 θ_0, θ_1



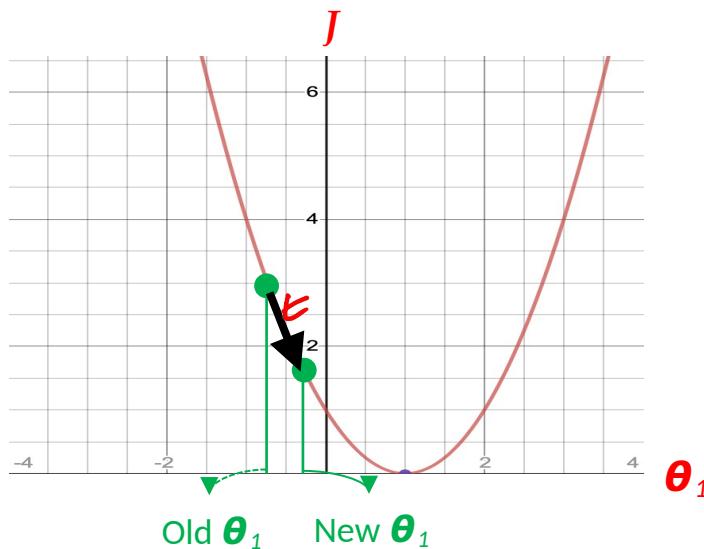
$\text{new } \theta_1 = \frac{\text{old}}{\theta_1 - \alpha} \frac{d J(\theta_1)}{d \theta_j}$

$= \theta_1 - \alpha \text{ (Negative Number)}$

Increase θ_1 by a certain value

The Impact of Partial Derivative

- optimization objective is to minimize $J(\theta_1)$
 θ_0, θ_1

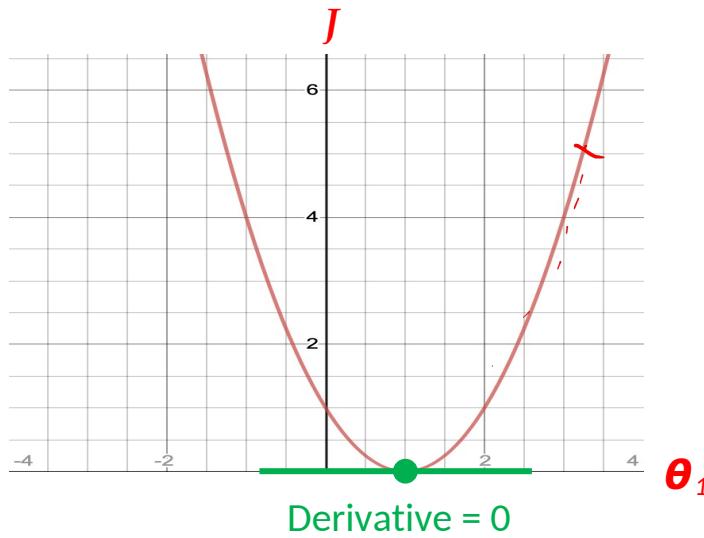


$$\begin{aligned}\theta_1 &= \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j} \\ &= \theta_1 - \alpha (\text{Negative Number})\end{aligned}$$

Increase θ_1 by a certain value

The Impact of Partial Derivative

- optimization objective is to minimize $J(\theta_1)$
 θ_0, θ_1



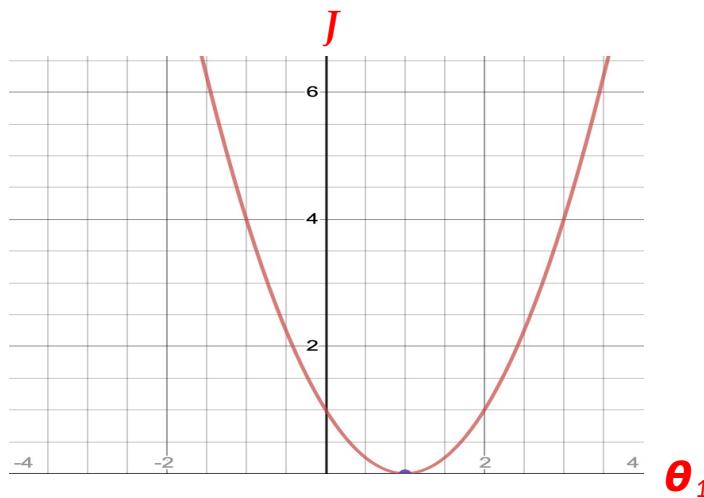
$$\begin{aligned}
 \text{new } \theta_1 &= \frac{\text{old } \theta_1 - \alpha}{d J(\theta_1)} \frac{d J(\theta_1)}{d \theta_j} \\
 &= \theta_1 - \alpha (\text{Zero})
 \end{aligned}$$

θ_1 remains the same, hence,
gradient descent converges

The Impact of Learning Rate

- optimization objective is to minimize $J(\theta_1)$
 θ_0, θ_1

Learning rate
hyper parameters



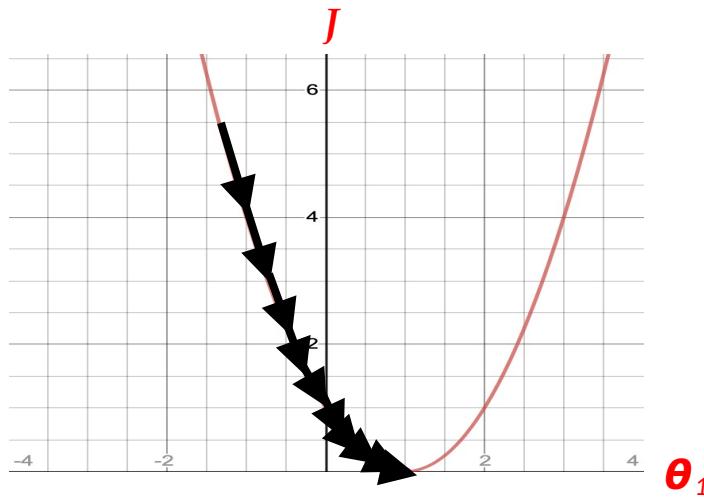
$$\theta_1 = \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j}$$

Learing Rate

What happens if α is too small?

The Impact of Learning Rate

- optimization objective is to minimize $J(\theta_1)$
 θ_0, θ_1



new ~~old~~ $\theta_1 = \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j}$
 $= \theta_1 - (\text{Too Small Number}) \frac{d J(\theta_1)}{d \theta_j}$

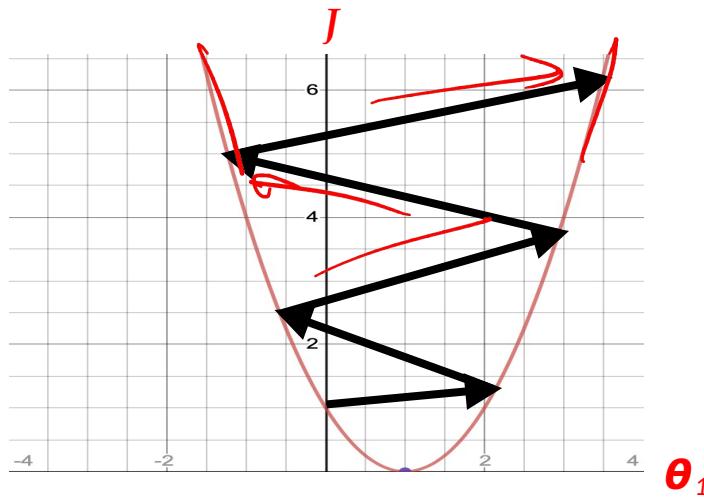
θ_1 changes only a tiny bit on each step,
 hence, gradient descent will render
 slow (will take more time to converge)

The Impact of Learning Rate

- optimization objective is to minimize $J(\theta_1)$
 θ_0, θ_1

$$\alpha = 0.1$$

$$0.01$$



$$\begin{aligned}
 \theta_1 &= \theta_1 - \alpha \frac{d J(\theta_1)}{d \theta_j} \\
 &= \theta_1 - (\text{Too Large Number}) \frac{d J(\theta_1)}{d \theta_j}
 \end{aligned}$$

θ_1 changes a lot (and probably faster) on each step, hence, gradient descent will potentially overshoot the minimum and, accordingly, fail to converge (or even diverge)

Gradient Descent algorithm :

Effect of Learning Rate

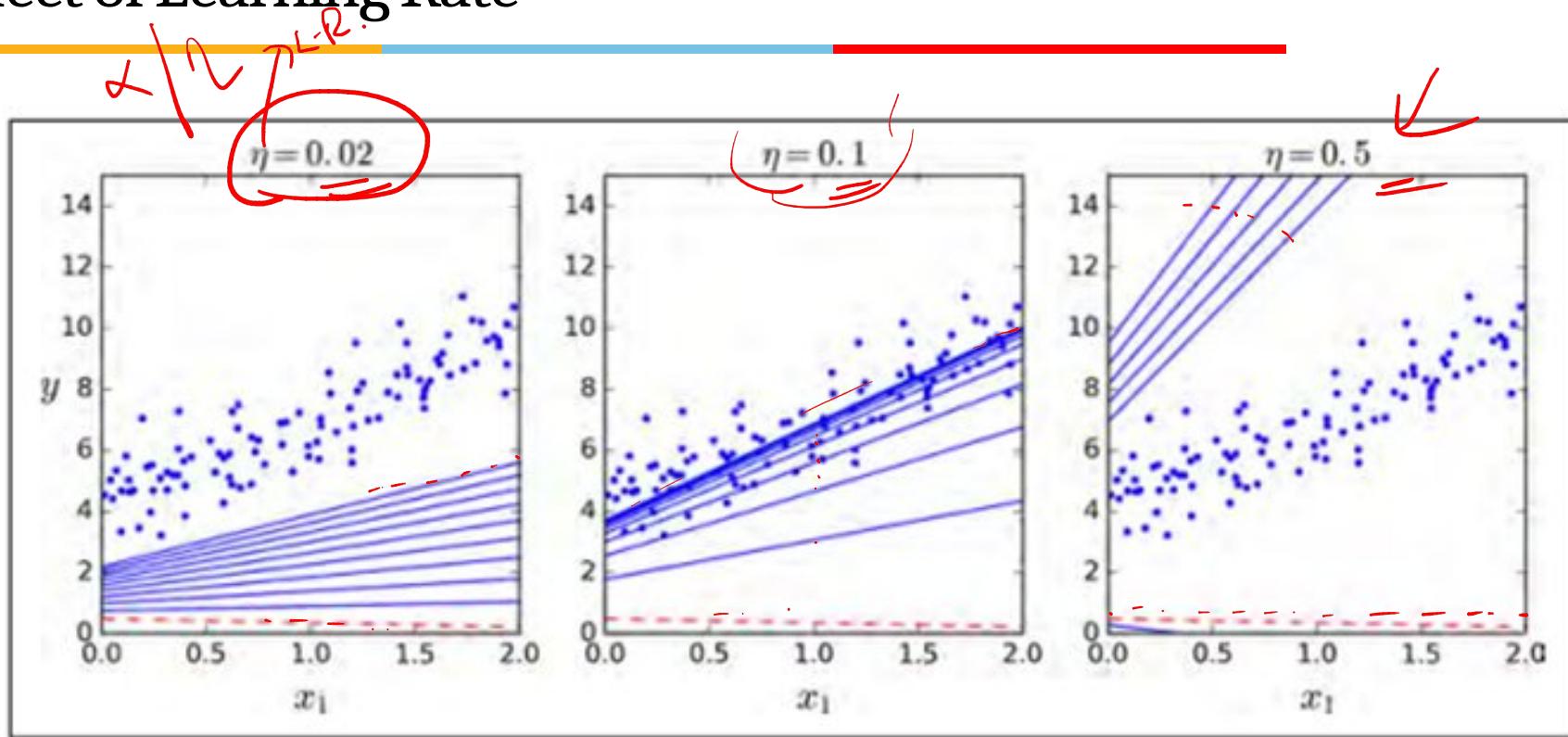


Figure 4-8. Gradient Descent with various learning rates

Gradient Descent algorithm :

Effect of Feature Scaling

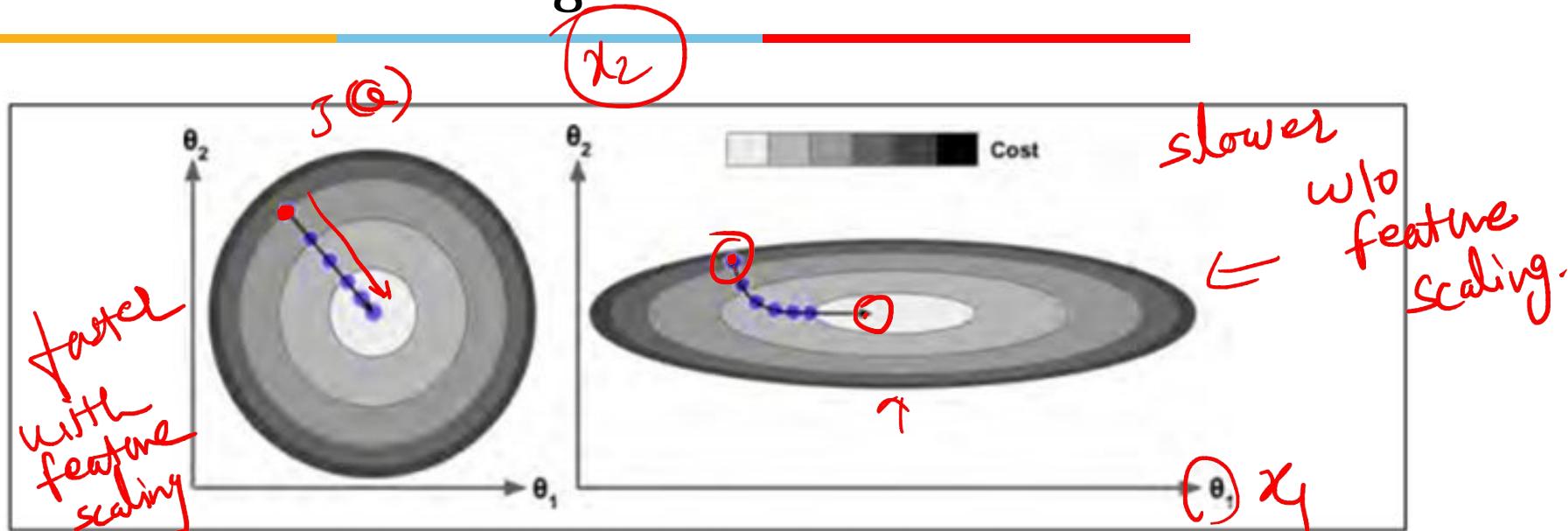


Figure 4-7. Gradient Descent with and without feature scaling

$$\nabla J(\theta) = \begin{pmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{pmatrix} = \begin{pmatrix} \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \end{pmatrix}$$

x_1 bedroom	x_2 size in sqft
900	500
1000	600
1100	700
1200	800
1300	900
1400	1000

Cost Function Linear Regression

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

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

mse

For insight on $J()$, let's assume $x \in \mathbb{R}$ so $\theta = [\theta_0, \theta_1]$

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

training example

Derivative of cost function for one training example

$$h(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x,$$

$$\begin{aligned}
 \frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_\theta(x) - y)^2 \\
 &= 2 \cdot \frac{1}{2} (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_\theta(x) - y) \\
 &= (h_\theta(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left(\sum_{i=0}^n \theta_i x_i - y \right) \\
 \frac{\partial J(\theta)}{\partial \theta_j} &= (h_\theta(x) - y) x_j = \frac{\partial}{\partial \theta_0} (\theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n - y)
 \end{aligned}$$

$x_0 \frac{\partial \theta_0}{\partial \theta_0} \Rightarrow x_0$
 $\theta_j \Rightarrow x_j$

Derivative of cost function for n training example

- Initialize θ
- Repeat until convergence

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

$x_j^{(i)}$

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

$$= \frac{1}{n} \sum_{i=1}^n \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right)$$

$$= \frac{1}{n} \sum_{i=1}^n \left(\sum_{k=0}^d \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}$$

θ_0

θ_1

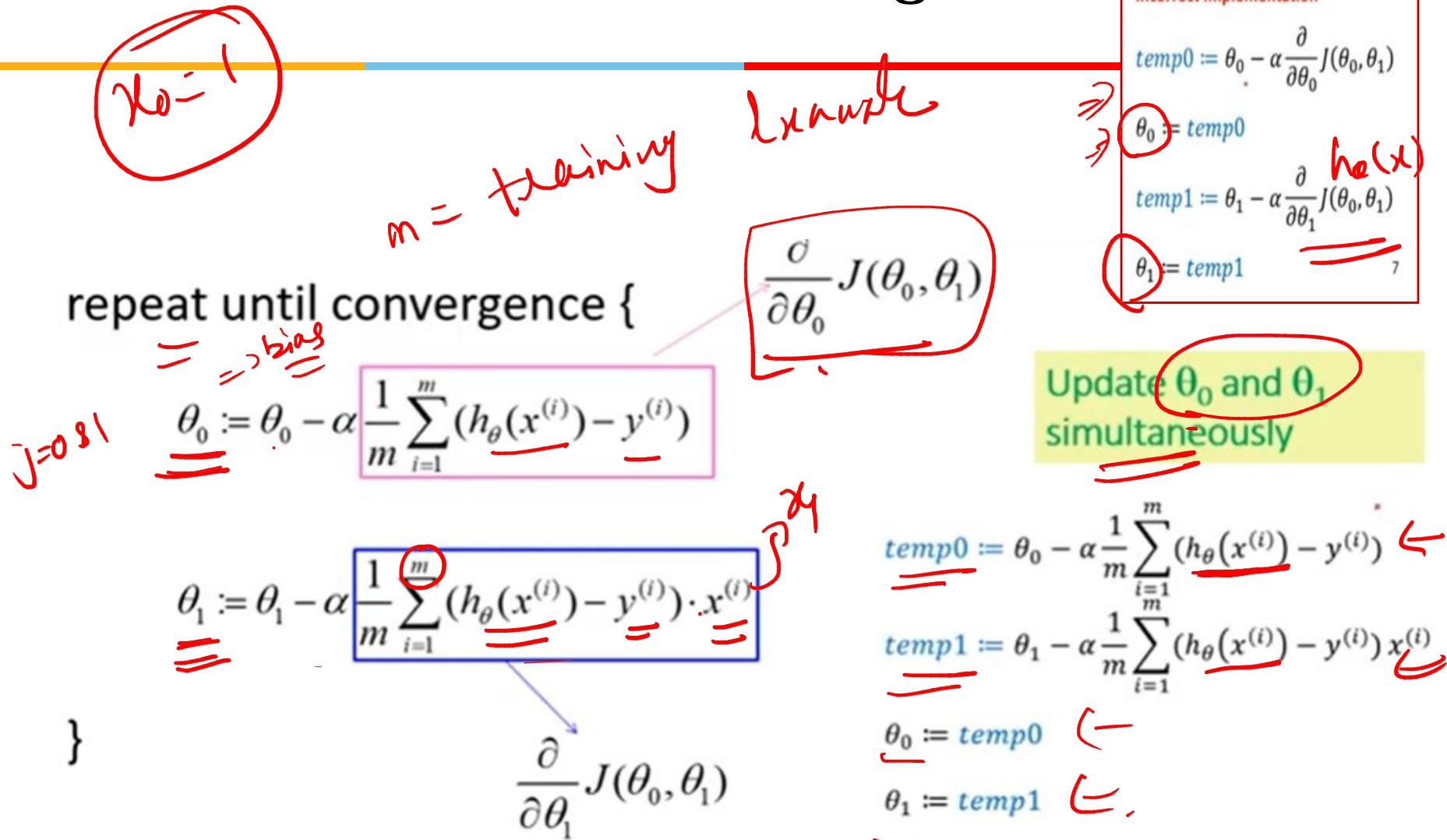
x_0

x_1

$d = \dim$

$n = \text{training examples}$

Gradient Descent for Linear Regression



Gradient Descent for Linear Regression

θ_j

$j=0 \text{ to } d$

θ

- Initialize θ

- Repeat until convergence

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

simultaneous update for $j = 0 \dots d$

(Pred - act) feature value

- 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 $\|\theta_{new} - \theta_{old}\|_2 < \epsilon$

$$L_2 \text{ norm: } \|\mathbf{v}\|_2 = \sqrt{\sum_i v_i^2} = \sqrt{v_1^2 + v_2^2 + \dots + v_{|v|}^2}$$

Gradient Descent: Variants

- **Batch** gradient descent refers to calculating the derivative from all training ~~data~~ before calculating an update.

Initialize the Parameters $(\theta_0^1, \theta_1^1, \dots)$

K=1

Repeat until Convergence {

$$\theta_0^{k+1} = \theta_0^k - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

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

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

.....

K=k+1

}

return $(\theta_0^1, \theta_1^1, \dots)$



Gradient Descent: Variants

- Minibatch refers to calculating derivative of mini groups of training data before calculating an update.

Divide the training instances into "N" batches each of size "m"

Initialize the Parameters ($\theta_0^1, \theta_1^1, \dots$)

K=1

Repeat until Convergence {

Repeat for every batch in 1 : N , each with 'm' instances {

$$\theta_0^{k+1} = \theta_0^k - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

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

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

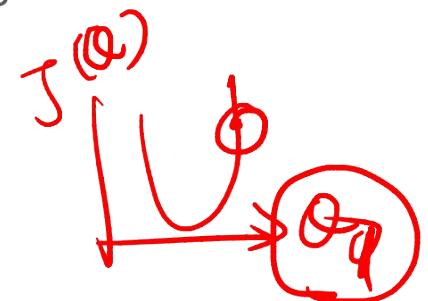
.....

K=k+1

}

}

return ($\theta_0^1, \theta_1^1, \dots$)



Gradient Descent: Variants

- **Stochastic gradient** descent refers to calculating the derivative from each training data instance and calculating the update immediately

Randomly shuffle training instances

Initialize the Parameters ($\theta_0^1, \theta_1^1, \dots$)

K=1

Repeat until Convergence {

Sample with replacement, only one random training instance “i” at a time

$$\left\{ \begin{array}{l} \theta_0^{k+1} = \theta_0^k - \alpha(h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1^{k+1} = \theta_1^k - \alpha(h_\theta(x^{(i)}) - y^{(i)}) * x_1^{(i)} \\ \theta_2^{k+1} = \theta_2^k - \alpha(h_\theta(x^{(i)}) - y^{(i)}) * x_2^{(i)} \\ \dots \end{array} \right.$$

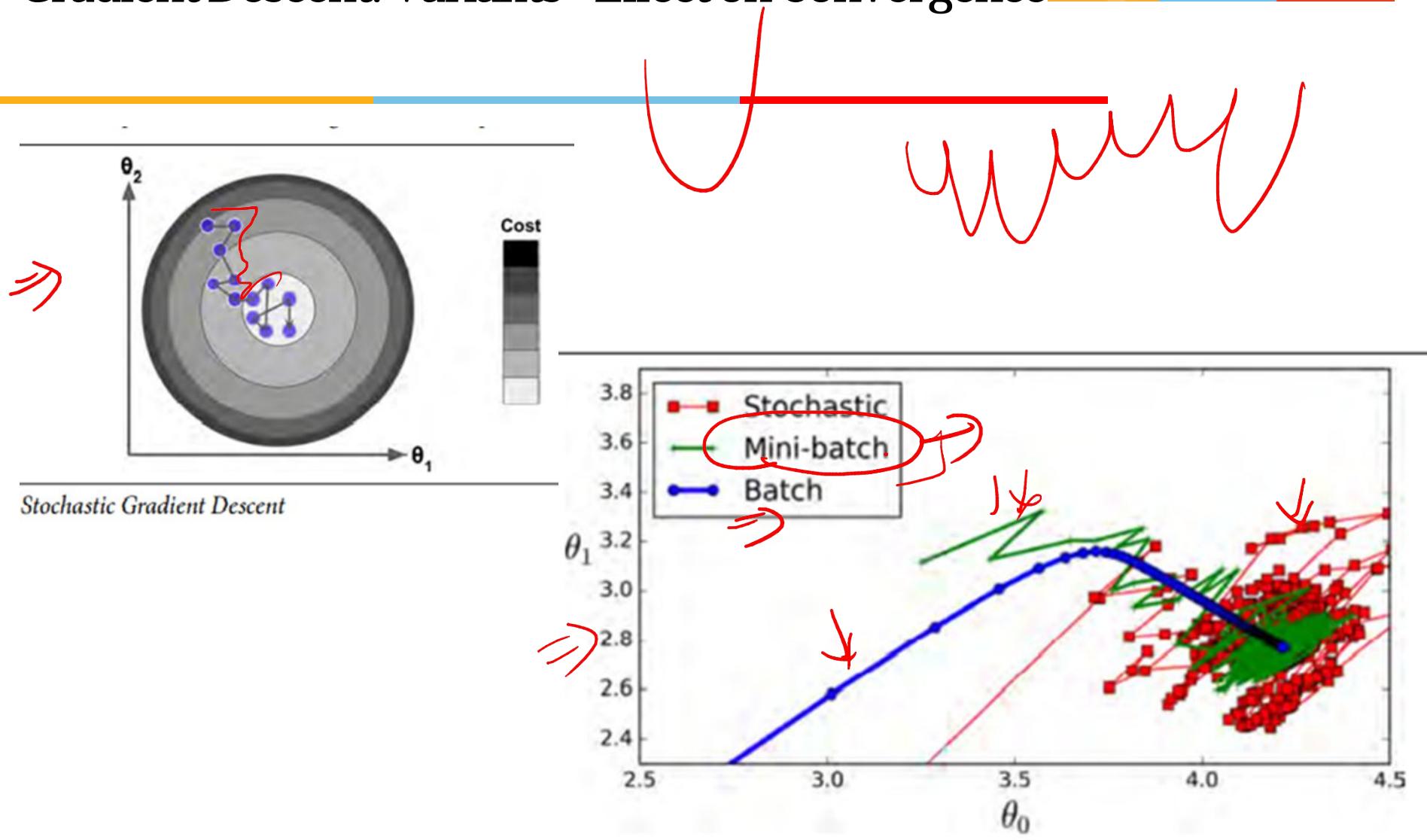
K=k+1

}

}

return ($\theta_0^1, \theta_1^1, \dots$)

Gradient Descent: Variants - Effect on Convergence



1. Gradient Descent paths in parameter space

Fit a linear Regression Line : Gradient Descent



Steps :

(Assuming : 'n' no.of.instances and two predictors $\{x_1, x_2\}$ and linear regression)

1. Identification of the equations $y = w_0 + w_1X_1 + W_2X_2$
2. Cost function & derivative
 1. $W_0` = w_0 - 1/n * \text{learning rate} * (\sum (w_0 + w_1X_1 + W_2X_2 - y))$
 2. $W_1` = w_1 - 1/n * \text{learning rate} * (\sum (w_0 + w_1X_1 + W_2X_2 - y) * x_1)$
 3. $W_2` = w_2 - 1/n * \text{learning rate} * (\sum (w_0 + w_1X_1 + W_2X_2 - y) * x_2)$
3. Apply the equations

Fit a linear Regression Line : Gradient Descent

Fit a linear regression. Show only the first iteration of Gradient descent algorithm using learning rate of **0.02** for the following data , if the Relative Risk of Coronary Heart Disease is believed to be only linearly dependent on BMI as well as Diastolic Pressure. Assume the intercept of the regression model as **5** and the slope of independent variables as **-0.03 (negative)**.

Patient	Systolic Pressure mm Hg	Diastolic Pressure mm Hg	BMI	Waist Circumference Threshold cm	RR-CHD (Relative Risk of Coronary Heart Disease)
1	140	80	35	100	1.81
2	120	80	25	80	1.22
3	130	100	30	60	1.71

Fit a linear Regression Line : Gradient Descent

- Initialize θ
- Repeat until convergence

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

simultaneous update
for j = 0 ... d

Patient	Systolic Pressure mm Hg	Diastolic Pressure mm Hg	BMI	Waist Circumference Threshold cm	RR-CHD (Relative Risk of Coronary Heart Disease)
1	140	80	35	100	1.81
2	120	80	25	80	1.22
3	130	100	30	60	1.71

Steps :

1. Identification of the equations: **RR-CHD = 5 - 0.03 * BMI - 0.03 * DiastolicPressure**
2. Cost function & derivative
 1. $W_0` = w_0 - 1/3 * 0.02 * (\text{sum}(5-0.03\text{BMI}-0.03\text{DiastolicPresure} - \text{RRCHD}))$
 $= 5 - 1/3 * 0.02 * (\text{sum}(5-0.03\text{BMI}-0.03\text{DiastolicPresure} - \text{RRCHD}))$
 2. $W_1` = w_1 - 1/3 * 0.02 * (\text{sum}(5-0.03\text{BMI}-0.03\text{DiastolicPresure} - \text{RRCHD}) * \text{BMI})$
 $= -0.03 - 1/3 * 0.02 * (\text{sum}(5-0.03\text{BMI}-0.03\text{DiastolicPresure} - \text{RRCHD}) * \text{BMI})$
 3. $W_2 = w_2 - 1/3 * 0.02 * (\text{sum}(5-0.03\text{BMI}-0.03\text{DiastolicPresure} - \text{RRCHD}) * \text{DiastolicPresure})$
 $= -0.03 - 1/3 * 0.02 * (\text{sum}(5-0.03\text{BMI}-0.03\text{DiastolicPresure} - \text{RRCHD}) * \text{DiastolicPresure})$
3. Apply the equations : Answer at **the end of first iteration:**
 $W_0 = 5.0016, W_1 = 0.0476, W_2 = 0.179$

References

- https://cs229.stanford.edu/lectures-spring2022/main_notes.pdf
- Pattern recognition and machine learning, Christopher bishop - CH3