

Université
Bretagne Sud

ubs:

Modèle de régression linéaire

Master 1

Traitement Numérique des Données

Sylvie Gibet

1

1

Housing Prices (Portland, OR)

Price
(in 1000s
of dollars)

Size (feet²)

Pour ce problème :

□ Apprentissage supervisé

□ donner la réponse juste pour chaque exemple contenu dans les données (valeur réelle, valeur discrète)

□ 1. Problème de régression

□ prédire la valeur réelle de sortie

□ 2. Classification

□ valeur de sortie discrète

2

2

Line Perret

1

Machine learning: deux étapes

- Deux étapes dans le processus de machine learning
 - ▣ Phase d'entraînement (training) : l'algorithme apprend sur une partie des données
 - ▣ Phase de test : test sur une autre partie des données

3

3

	Taille en m2	Prix en 1000 * €
Training set (apprentissage)	189	386
prix des maisons	127	106
	138	265
	76	149

...

- Objectif : apprendre à partir de ces données comment prédire les coûts des maisons
- Notation
 - ▣ m: nombre d'exemples d'apprentissage
 - ▣ x: variables d'entrée / features
 - ▣ y: variable de sortie/ target variable
 - (x, y) : un exemple d'apprentissage
 - ($x^{(i)}$, $y^{(i)}$) : ith exemple de l'apprentissage (ith ligne)

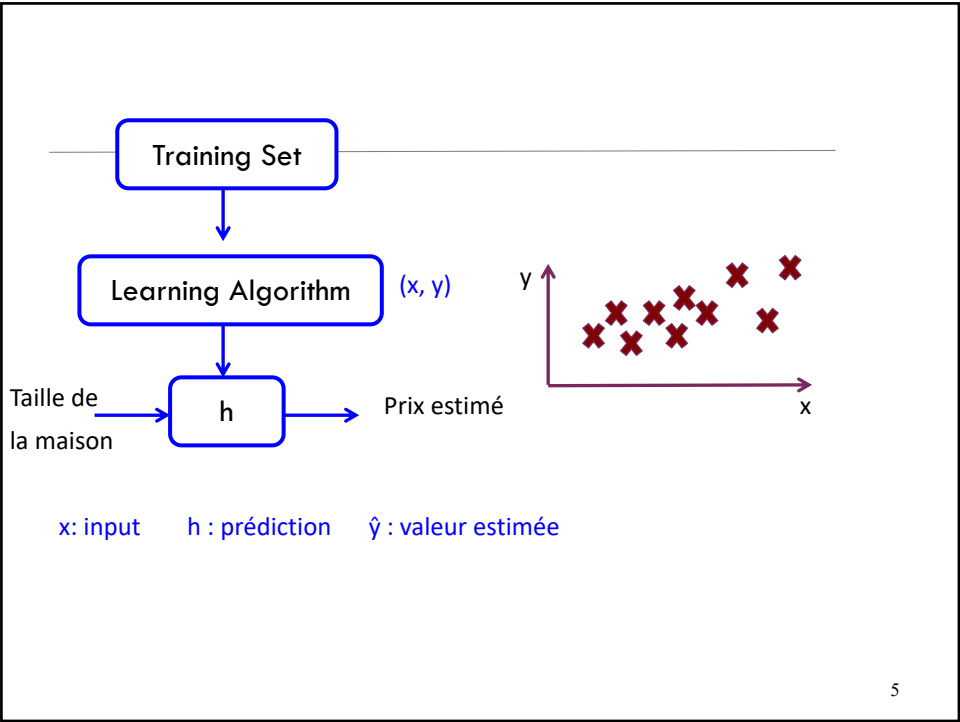
$x^{(1)} = 189$

$x^{(2)} = 127$

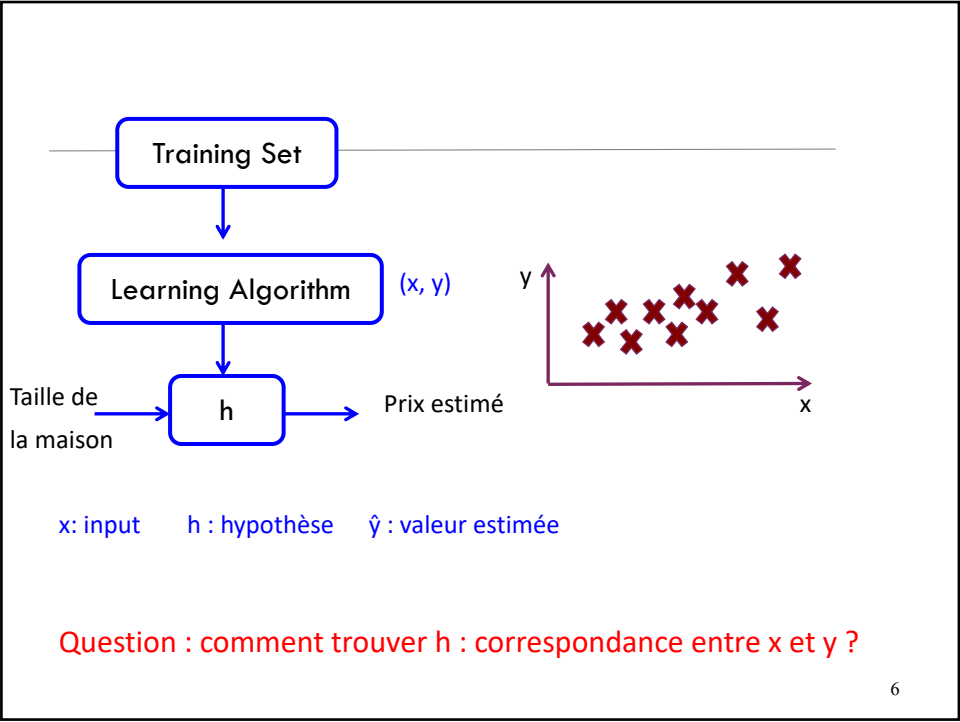
$y^{(1)} = 386$

4

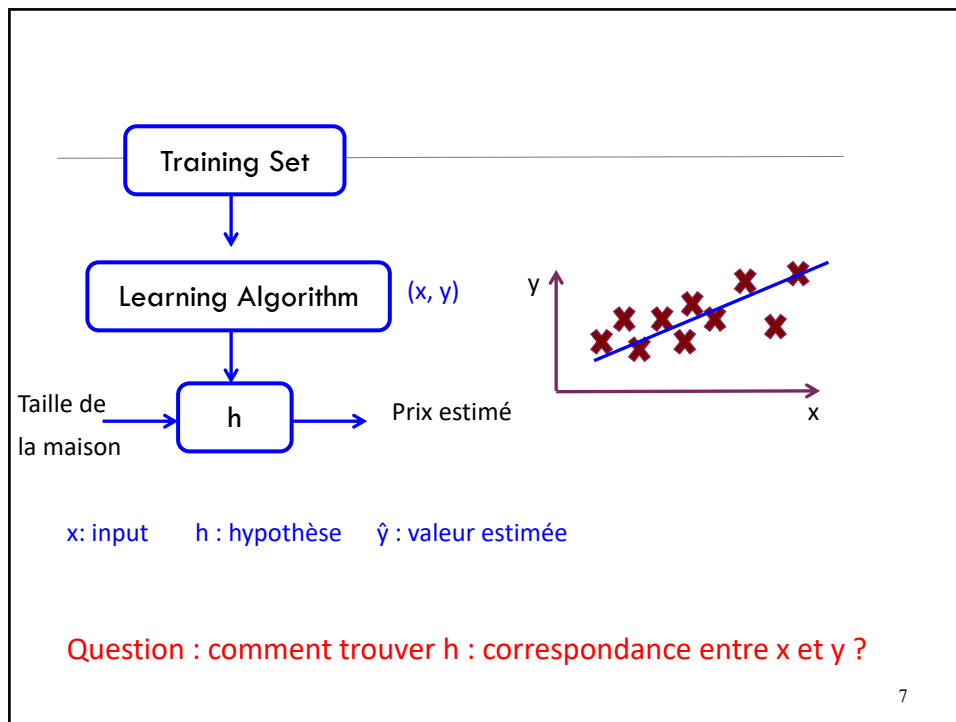
4



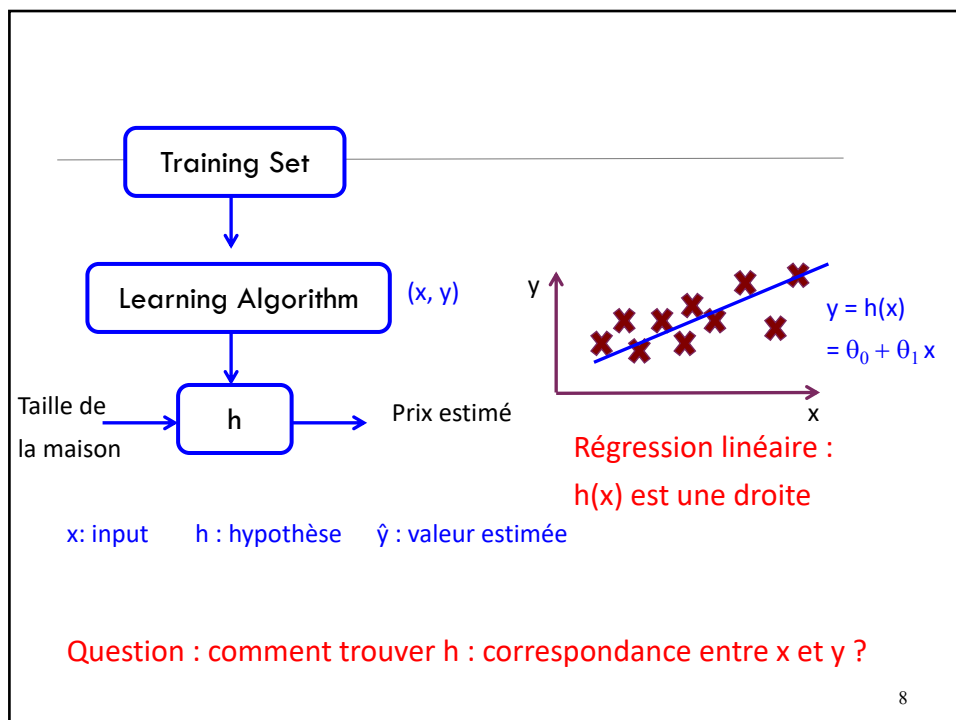
5



6



7



8

	Taille en m2	Prix en 1000 * €
Training set (apprentissage)	189	386
prix des maisons	127	106
	138	265
	76	149
	...	

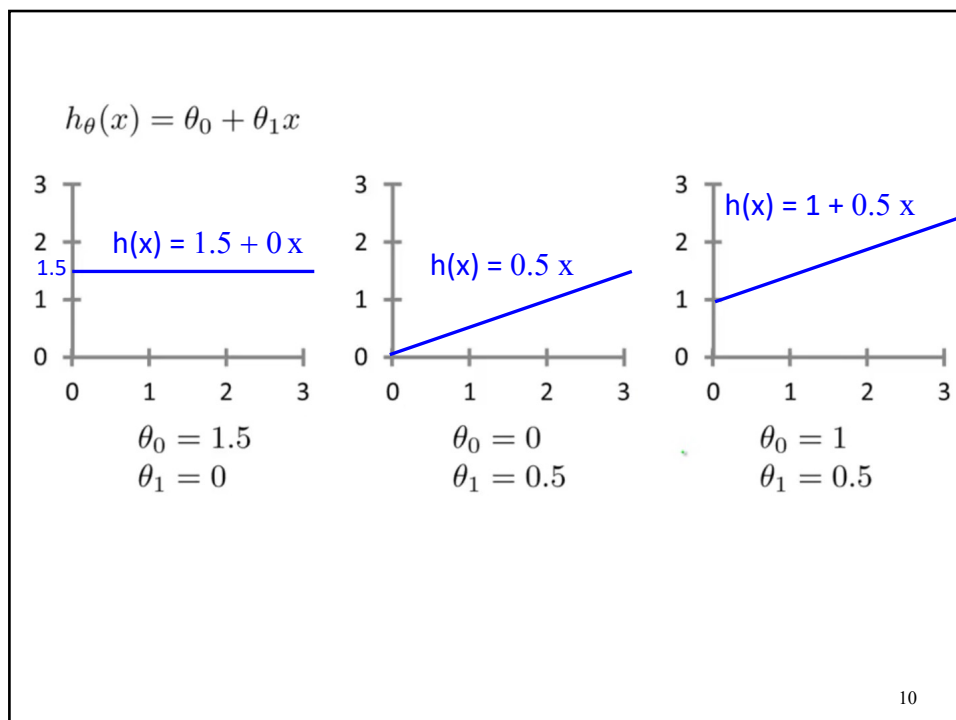
Comment choisir la fonction de mapping h ?

Hypothèse : on prend une droite (fonction linéaire de la variable x) :

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

x : variable d'entrée de la régression linéaire, ici la taille des maisons
 \hat{y} : variable estimée (sortie) : prix des maisons
 θ_i : paramètres (2 paramètres, θ_0 et θ_1)
 Comment choisir les paramètres θ_0 et θ_1 ?

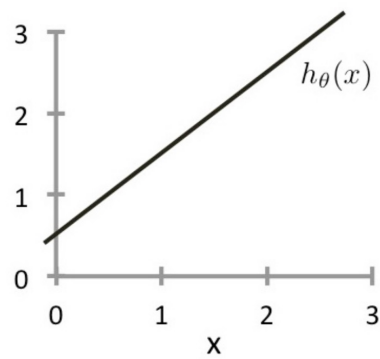
9



10

Exercice : Soit le diagramme ci-dessous : $h_{\theta}(x) = \theta_0 + \theta_1 x$.

Quelles sont les valeurs de θ_0 et θ_1 ?



- A. $\theta_0 = 0, \theta_1 = 1$
- B. $\theta_0 = 0.5, \theta_1 = 1$
- C. $\theta_0 = 1, \theta_1 = 0.5$
- D. $\theta_0 = 1, \theta_1 = 1$

11

11



Fonction coût

12

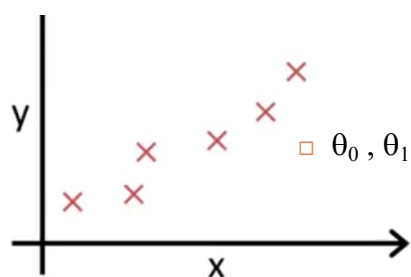
12

Fonction coût : définition

- Fonction qui va nous permettre d'évaluer la qualité de la droite qui approxime des données
- Peut s'exprimer comme une distance entre les données et la droite
 - ▣ Se ramène à la somme des distances entre chaque exemple contenu dans les données et sa "projection" sur la droite

13

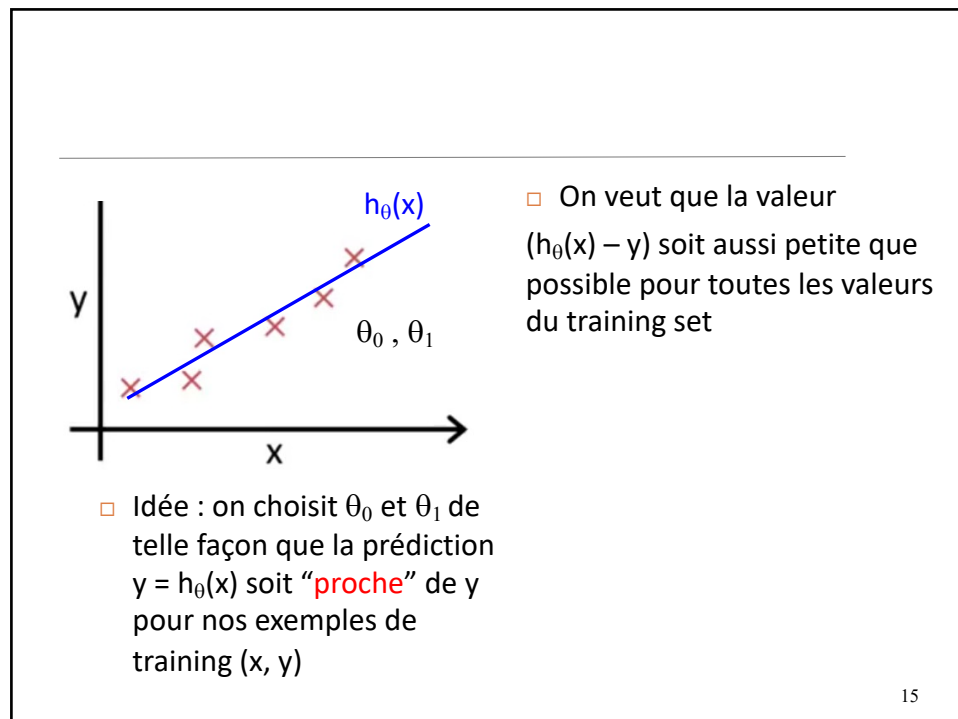
13



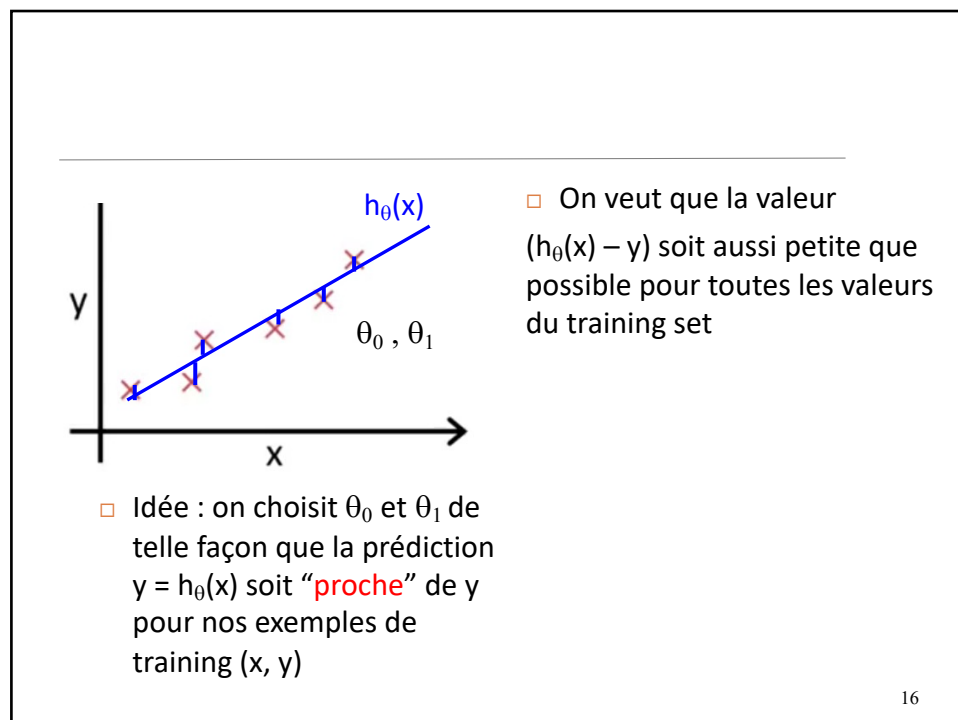
- Idée : on choisit θ_0 et θ_1 de telle façon que la prédiction $y = h_0(x)$ soit "proche" de y pour nos exemples de training (x, y)

14

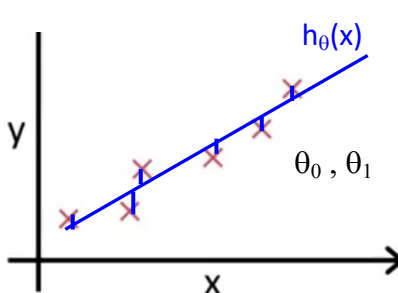
14



15



16



□ On veut que la valeur $(h_{\theta}(x) - y)$ soit aussi petite que possible pour toutes les valeurs du training set

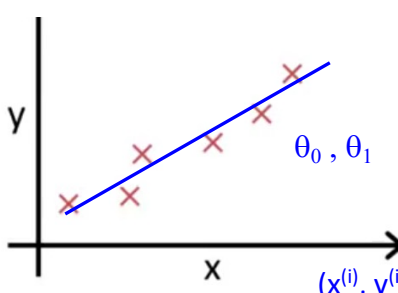
□ Une façon de faire est de minimiser le coût quadratique :

$(h_{\theta}(x) - y)^2$ pour tout le training set, suivant les paramètres θ_0, θ_1

□ Idée : on choisit θ_0 et θ_1 de telle façon que la prédiction $y = h_{\theta}(x)$ soit “proche” de y pour nos exemples de training (x, y)

17

17



□ On veut minimiser / (θ_0, θ_1) :

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

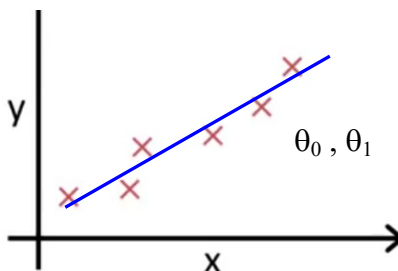
$y^{(i)} = h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$

$J(\theta_0, \theta_1)$: Fonction coût

□ Idée : on choisit θ_0 et θ_1 de telle façon que la prédiction $y = h_{\theta}(x)$ soit “proche” de y pour nos exemples de training (x, y)

18

18



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

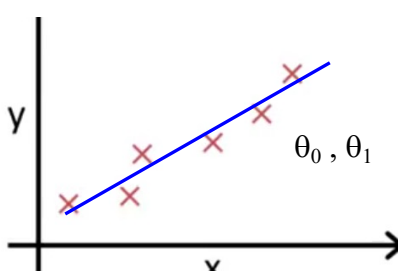
Minimiser $J(\theta_0, \theta_1)$ suivant θ_0, θ_1

$J(\theta_0, \theta_1)$: fonction coût encore appelée **erreur quadratique**

□ Idée : on choisit θ_0 et θ_1 de telle façon que la prédiction $y = h_{\theta}(x)$ soit “proche” de y pour nos exemples de training (x, y)

19

19



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

Minimiser $J(\theta_0, \theta_1)$ suivant θ_0, θ_1

$J(\theta_0, \theta_1)$: fonction coût encore appelée **erreur quadratique**

□ Idée : on choisit θ_0 et θ_1 de telle façon que la prédiction $y = h_{\theta}(x)$ soit “proche” de y pour nos exemples de training (x, y)

Bonne fonction d'erreur pour beaucoup de problèmes de régression

20

20

Fonction de coût Intuition I un seul paramètre θ_1

21

21

□ Hypothèse : $h_{\theta}(x) = \theta_0 + \theta_1 x$ avec
 x : input, $y = h_{\theta}(x)$: output

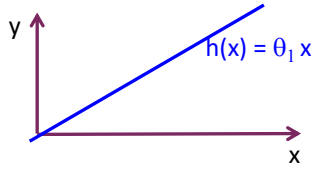
□ Paramètres : θ_0, θ_1

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

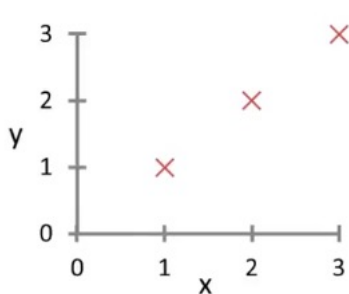
□ But : minimiser $J(\theta_0, \theta_1)$
 θ_0, θ_1

22

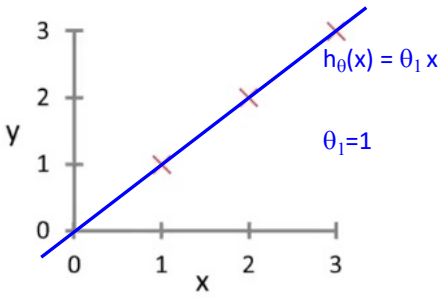
22

<div><div><div>□ Hypothèse : $y = h_{\theta}(x) = \theta_0 + \theta_1 x$</div><div>□ Paramètres : θ_0, θ_1</div><div>□ Fonction coût :</div><div>$J(\theta_0, \theta_1) = \frac{1}{2m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$</div><div>□ But : minimiser $J(\theta_0, \theta_1)$</div><div>θ_0, θ_1</div></div></div>	<div><div>Problème simplifié</div><div>□ $h_{\theta}(x) = \theta_1 x$</div><div>$\theta_0 = 0$</div><div></div><div>$J(\theta_1) = \frac{1}{2m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$</div><div>$\theta_1 x^{(i)}$</div><div>minimiser $J(\theta_1)$</div><div>θ_1</div></div>
--	---

23

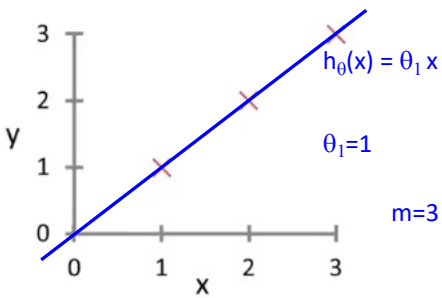
<div><div>$h_{\theta}(x)$</div><div>(pour θ_1 fixe, c'est une fonction de x)</div><div></div></div>	<div><div>$J(\theta_1)$</div><div>(fonction du paramètre θ_1)</div></div>
--	--

24

$h_{\theta}(x)$ (pour θ_1 fixe, c'est une fonction de x)	$J(\theta_1)$ (fonction du paramètre θ_1)
	

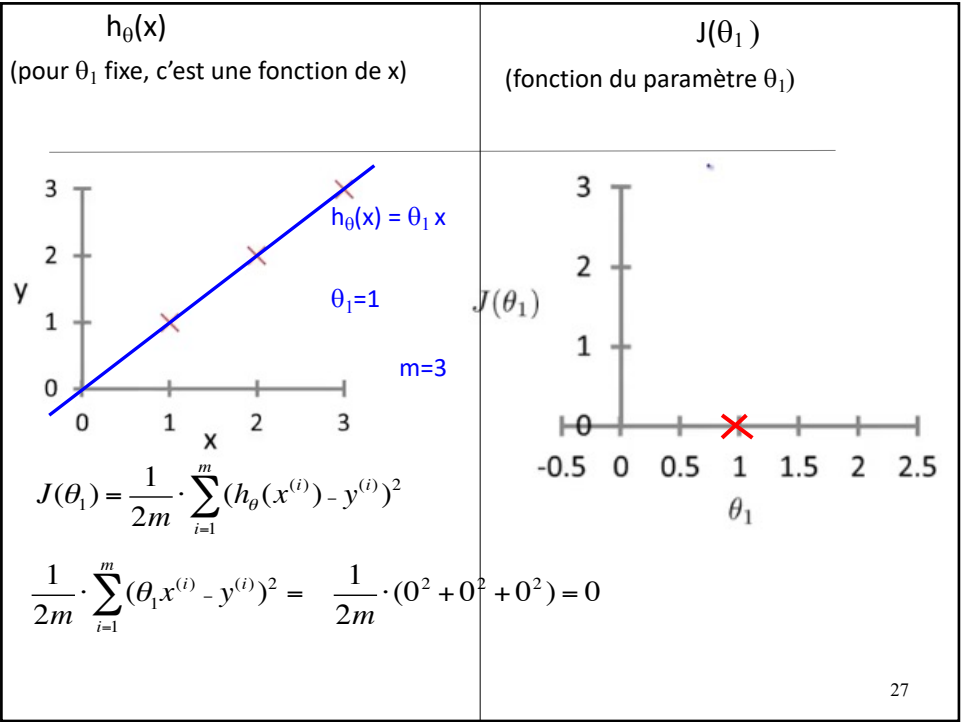
25

25

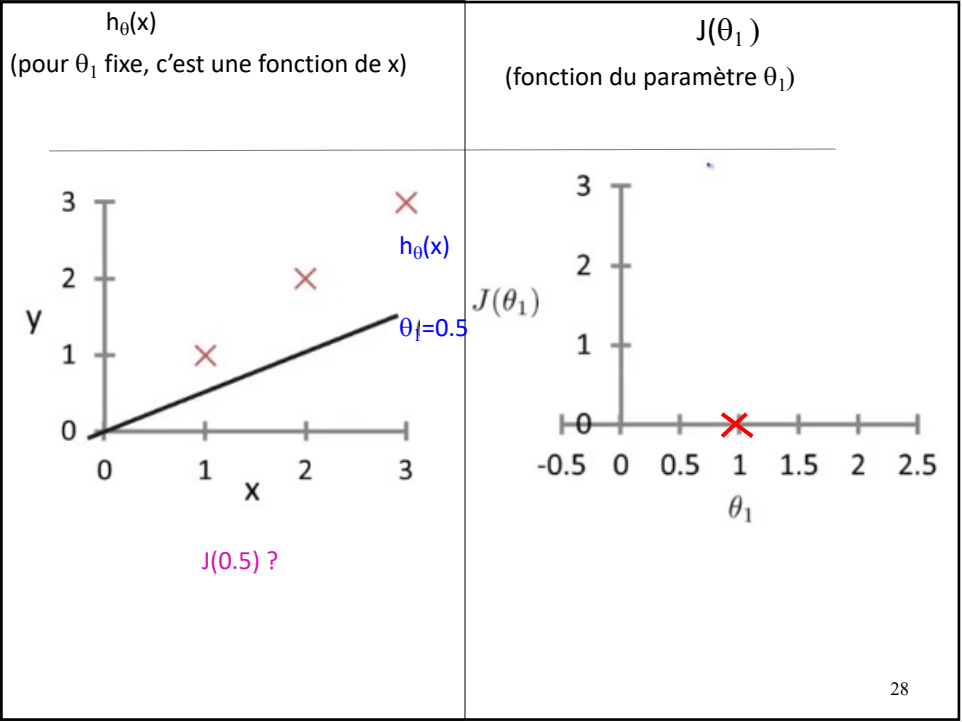
$h_{\theta}(x)$ (pour θ_1 fixe, c'est une fonction de x)	$J(\theta_1)$ (fonction du paramètre θ_1)
 <div>$J(\theta_1) = \frac{1}{2m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$\frac{1}{2m} \cdot \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 = \frac{1}{2m} \cdot (0^2 + 0^2 + 0^2) = 0$</div>	

26

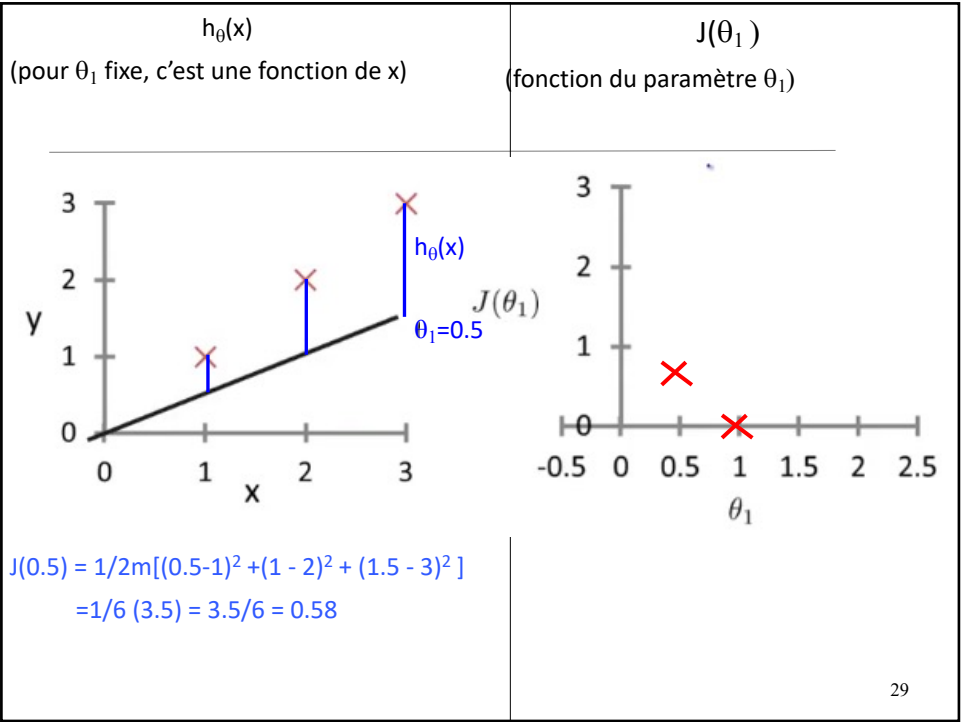
26



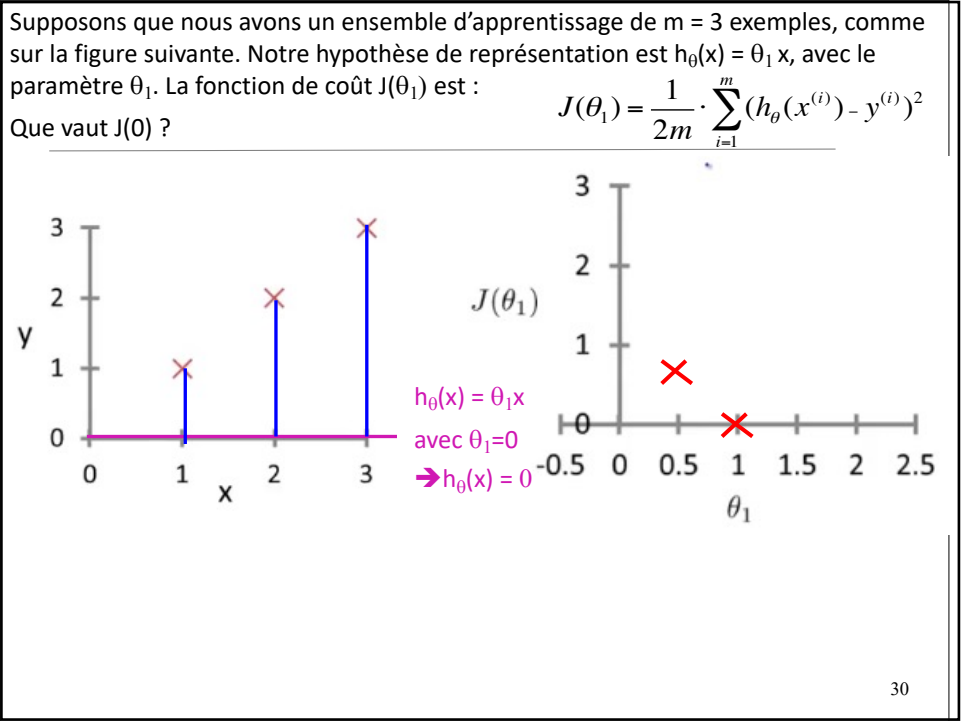
27



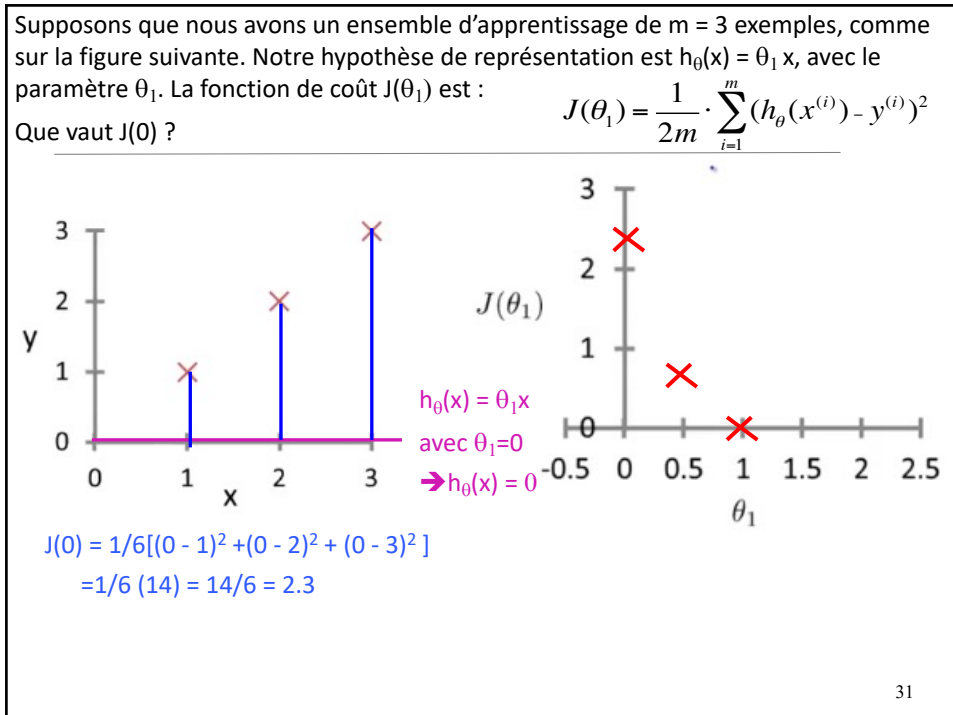
28



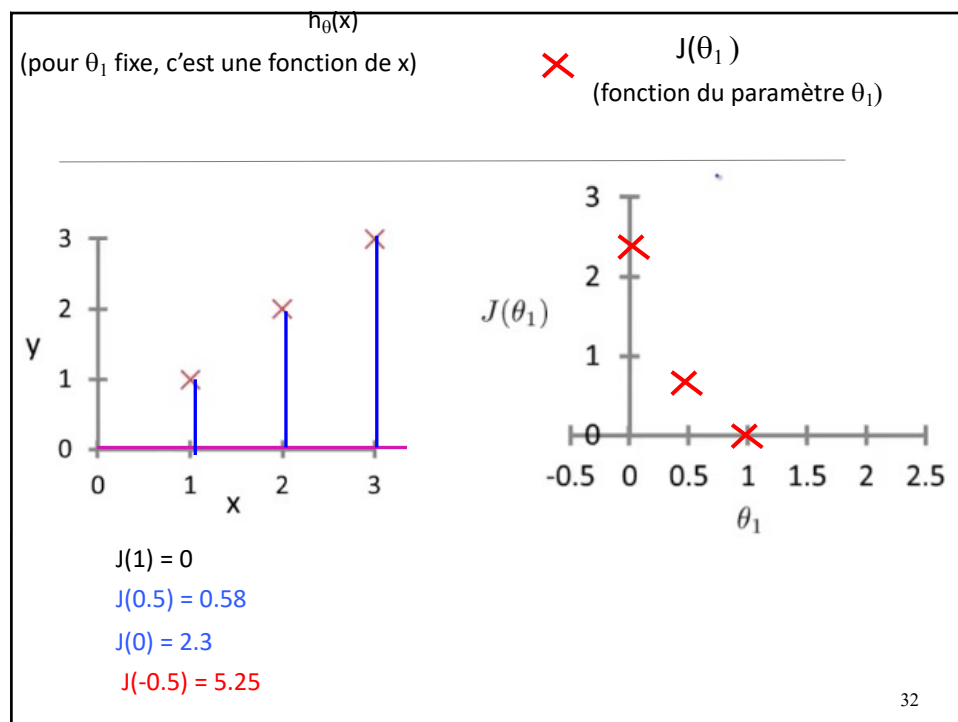
29



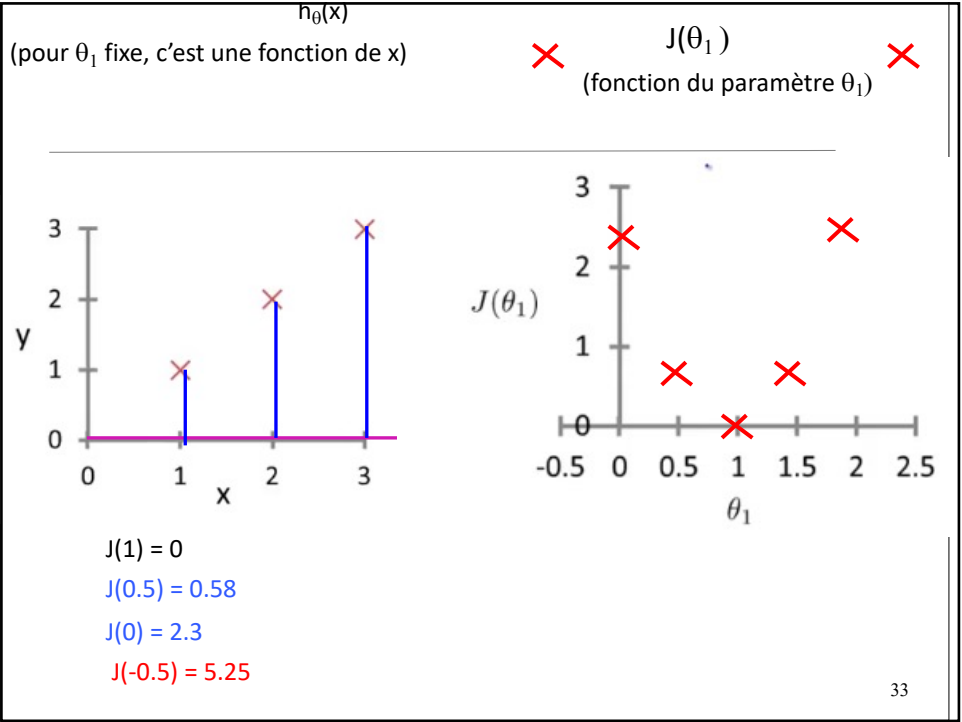
30



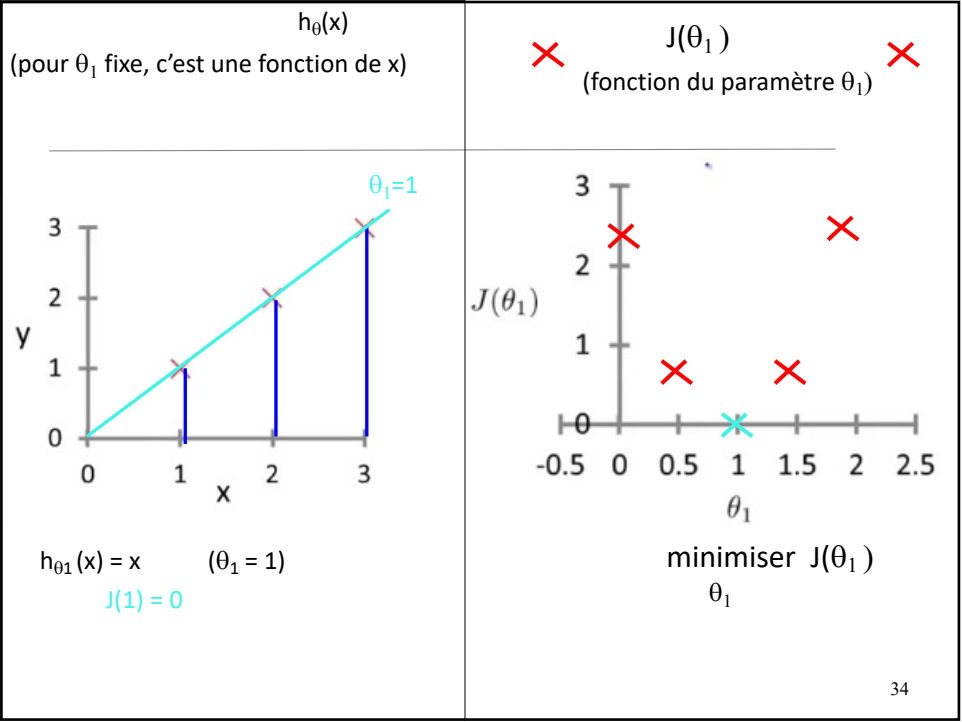
31



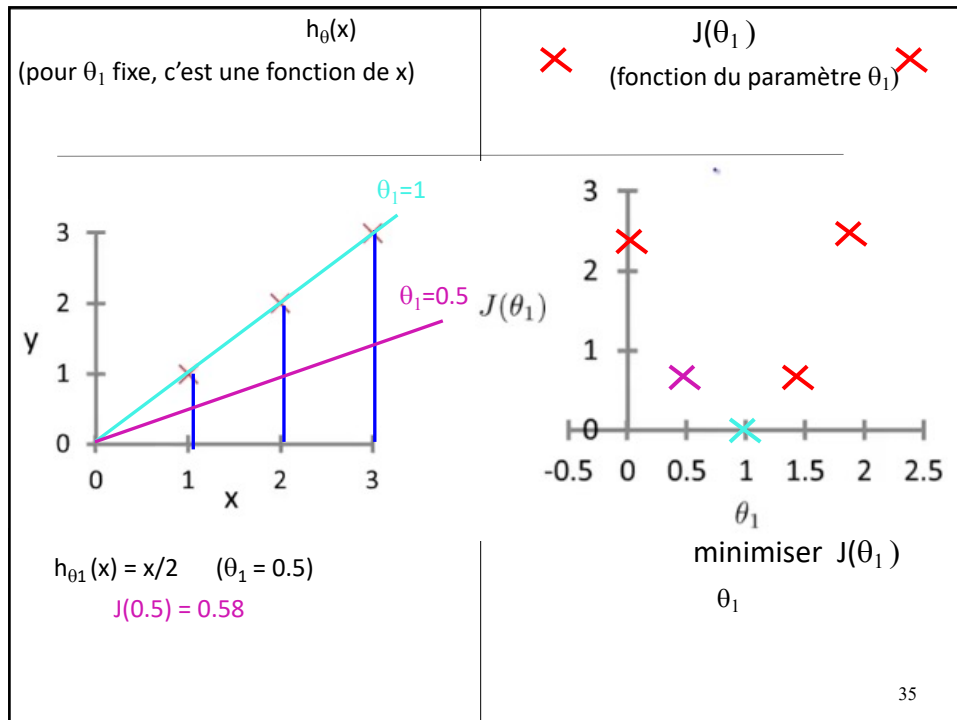
32



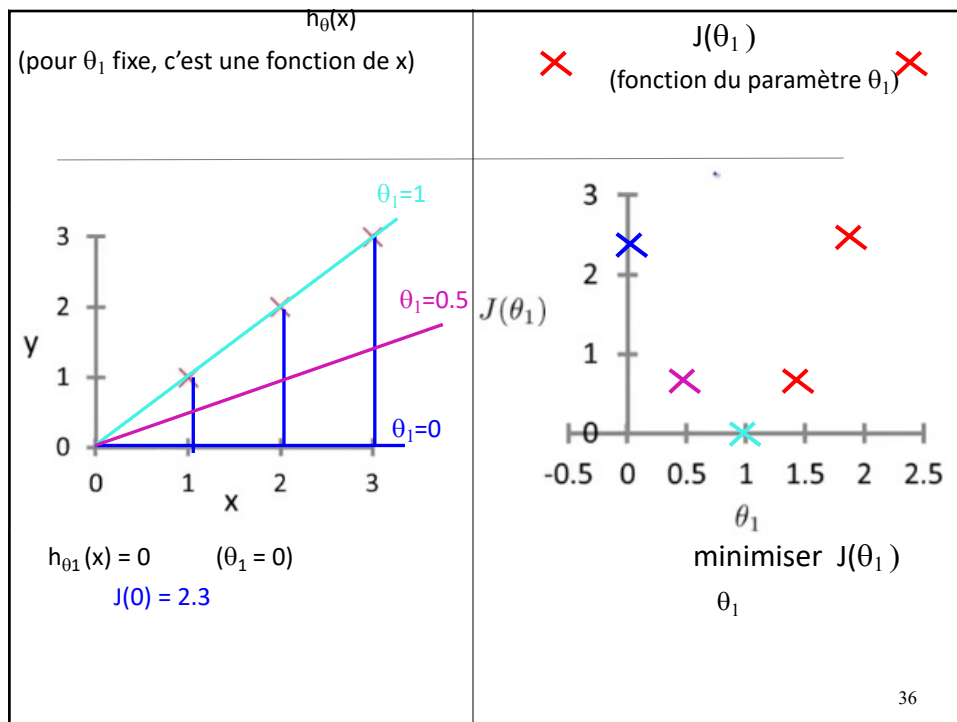
33



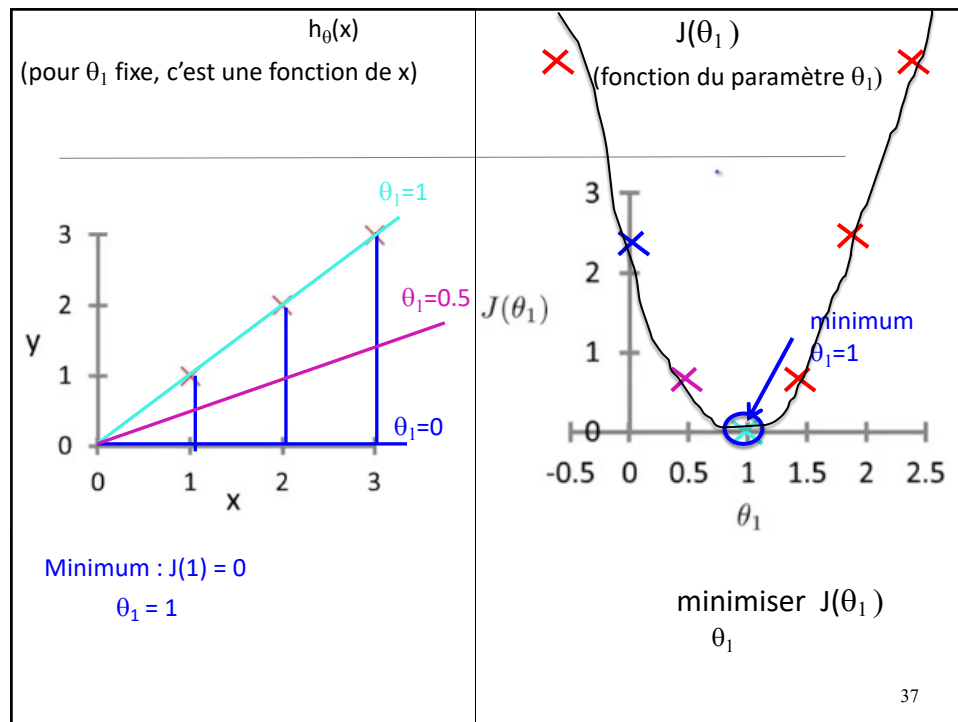
34



35



36



37

Fonction de coût

Intuition II

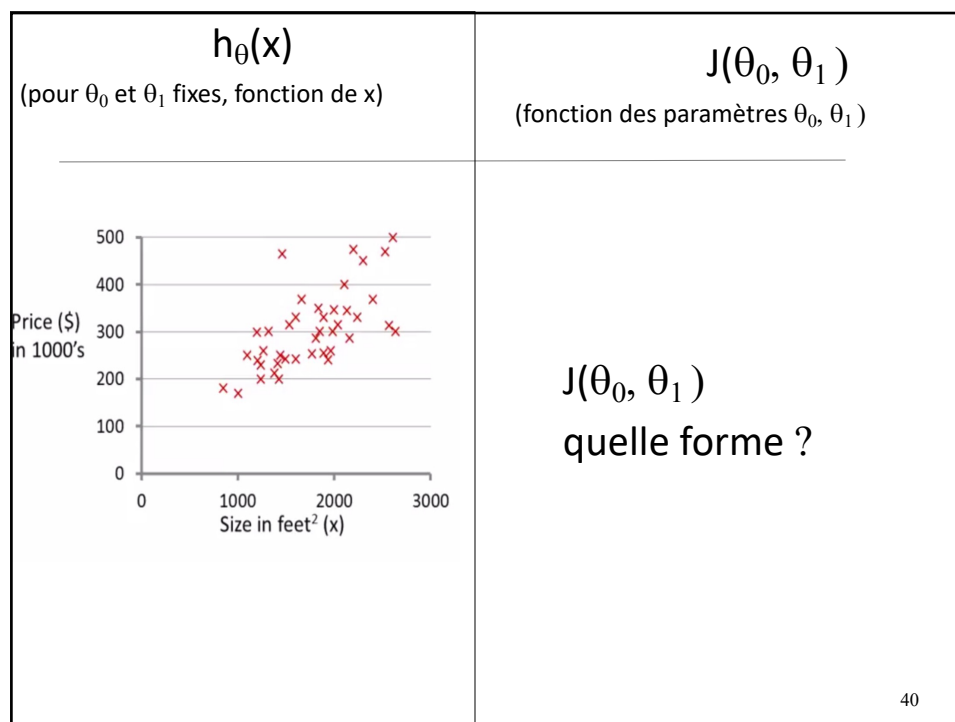
deux paramètres θ_0, θ_1

38

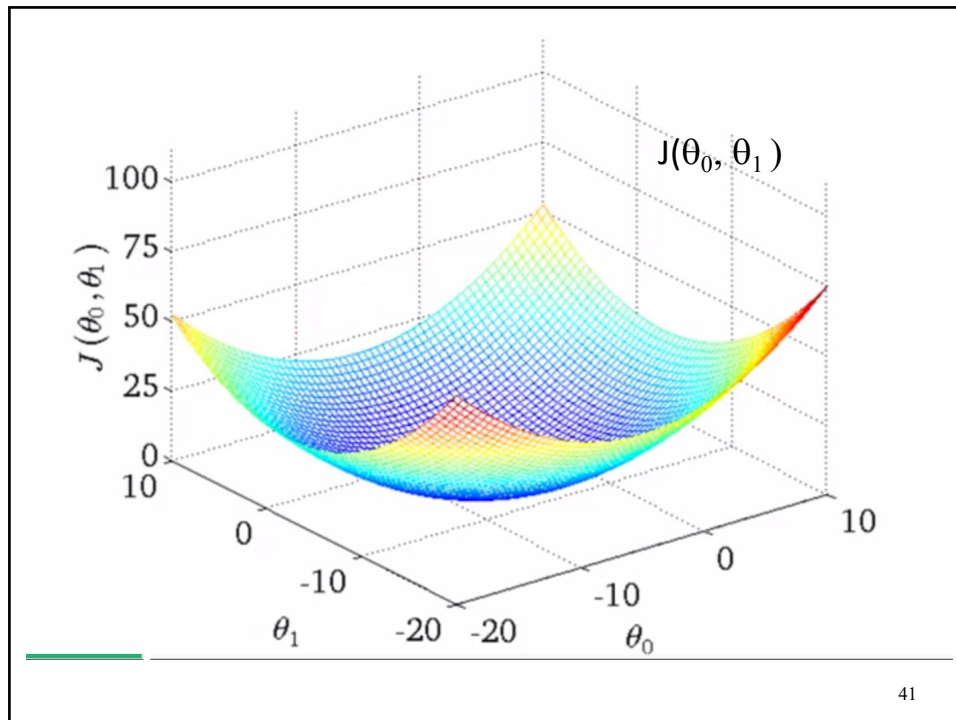
- Hypothèse : $h_{\theta}(x) = \theta_0 + \theta_1 x$
- Paramètres : θ_0, θ_1
- Fonction coût : $J(\theta_0, \theta_1) = \frac{1}{2m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- But : minimiser $J(\theta_0, \theta_1)$
 θ_0, θ_1

39

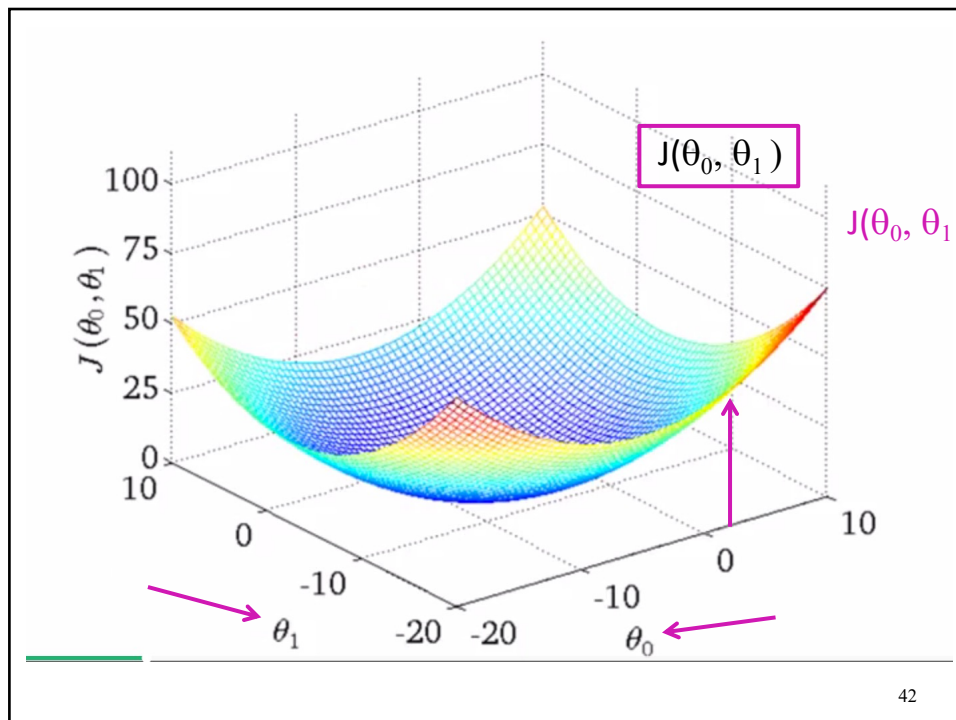
39



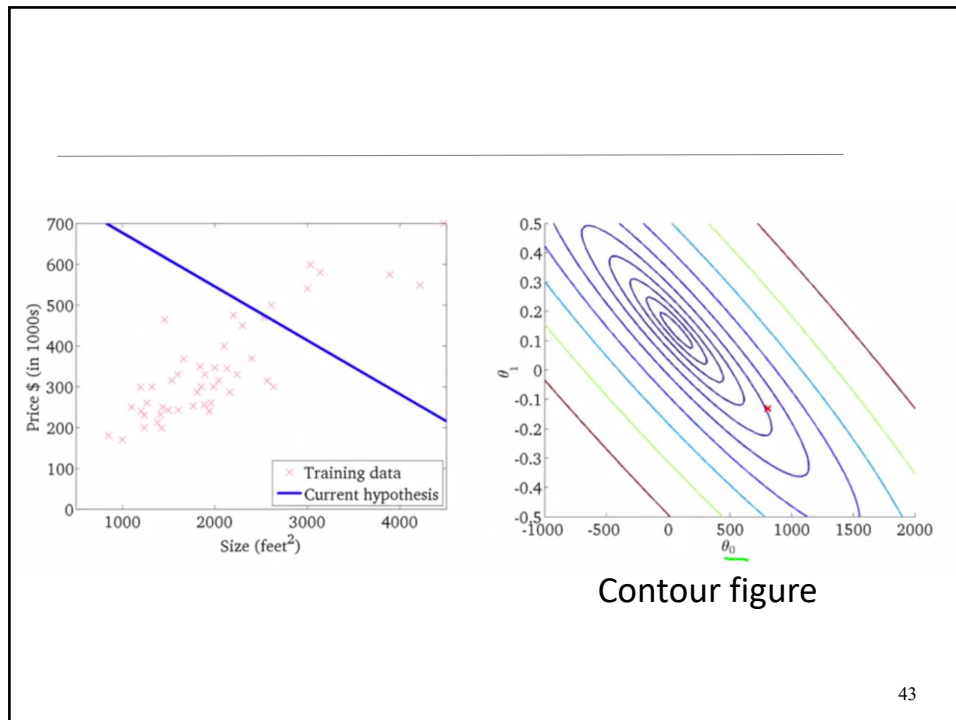
40



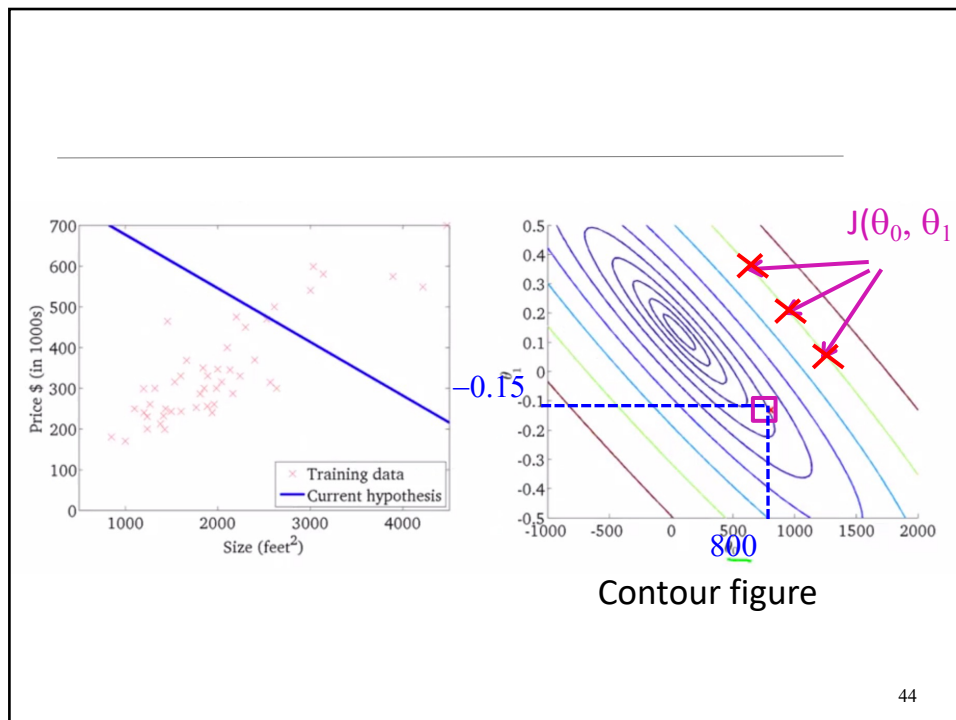
41



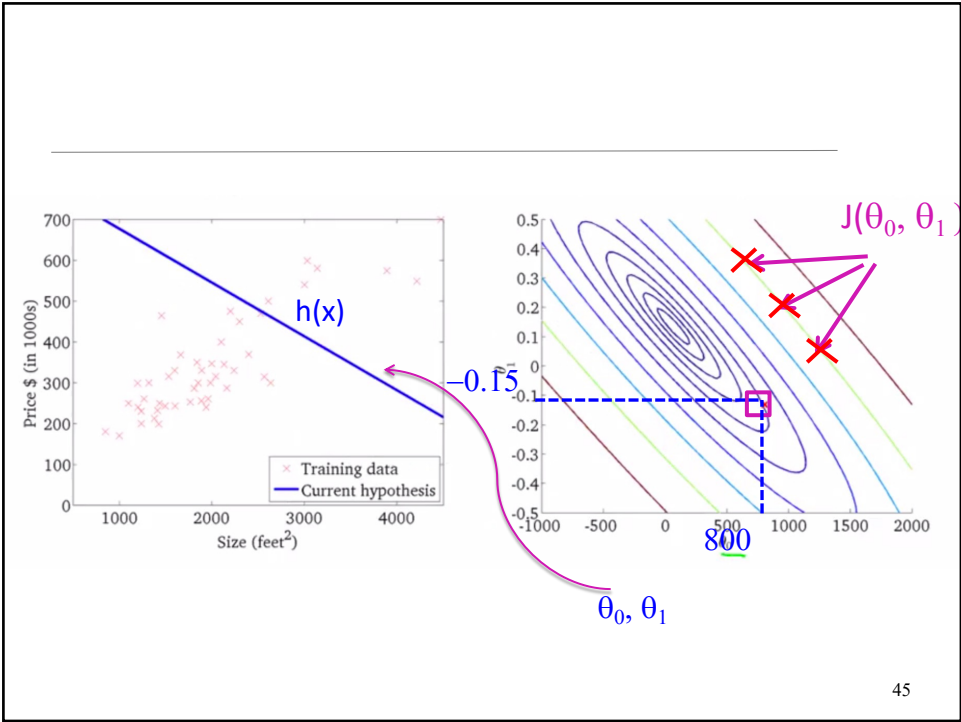
42



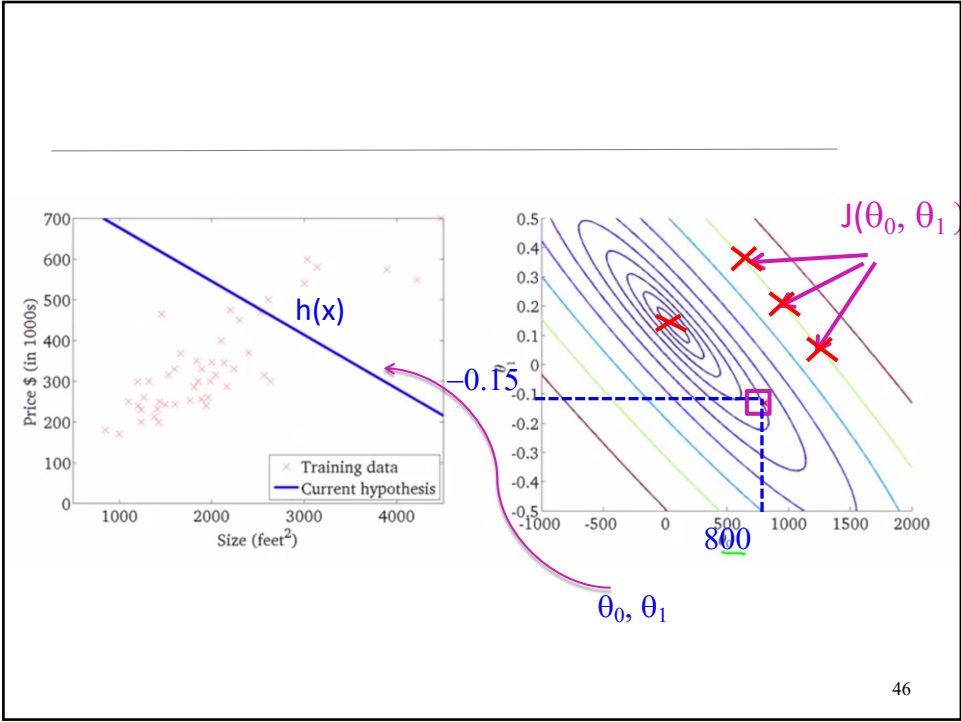
43



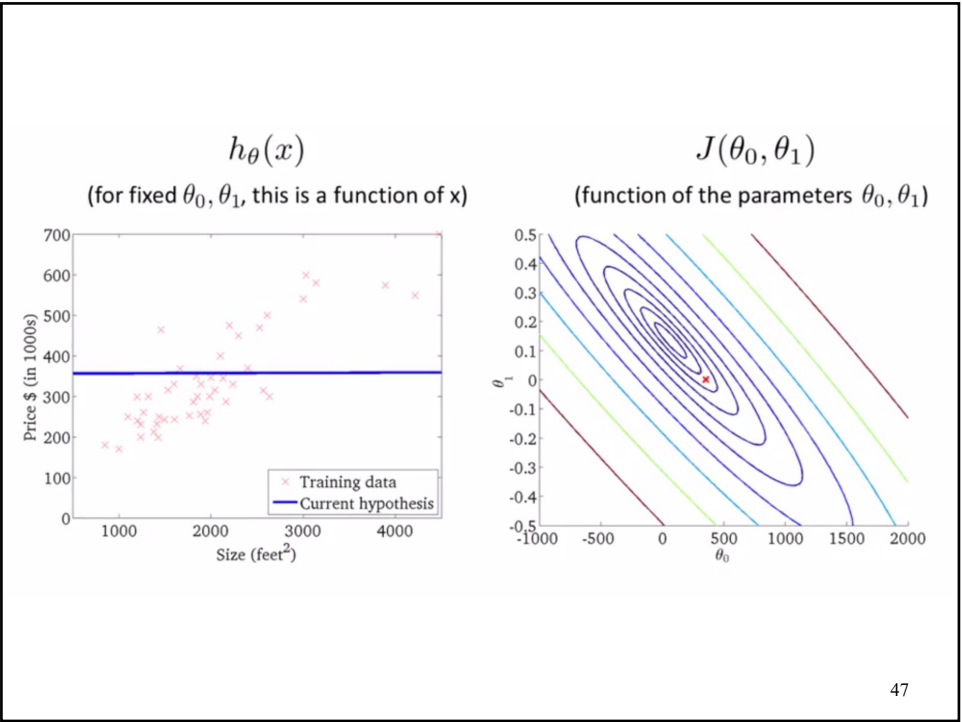
44



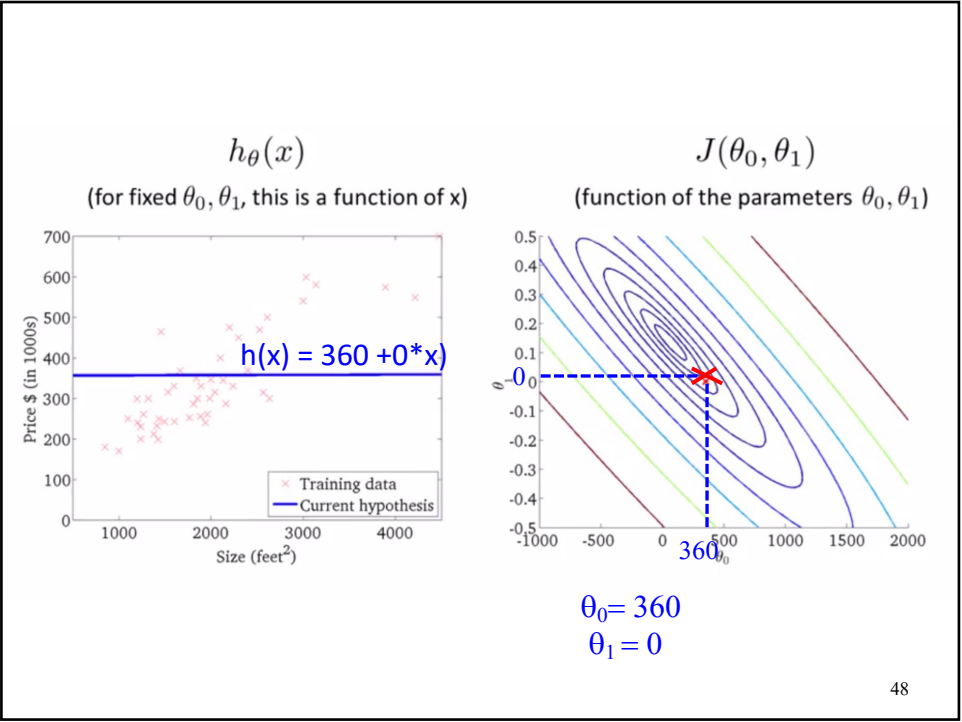
45



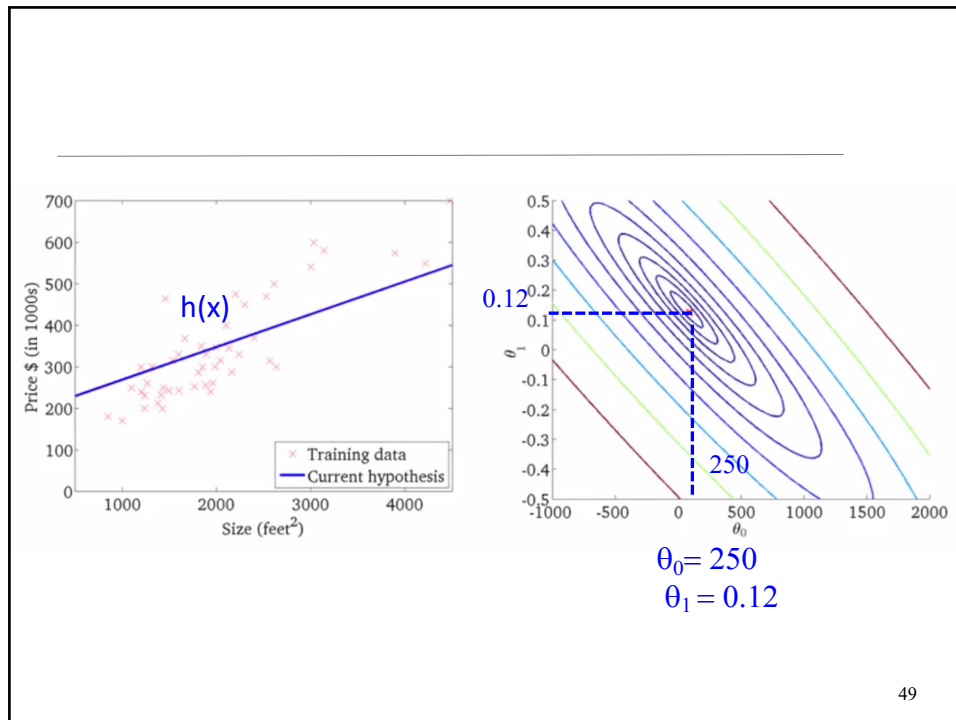
46



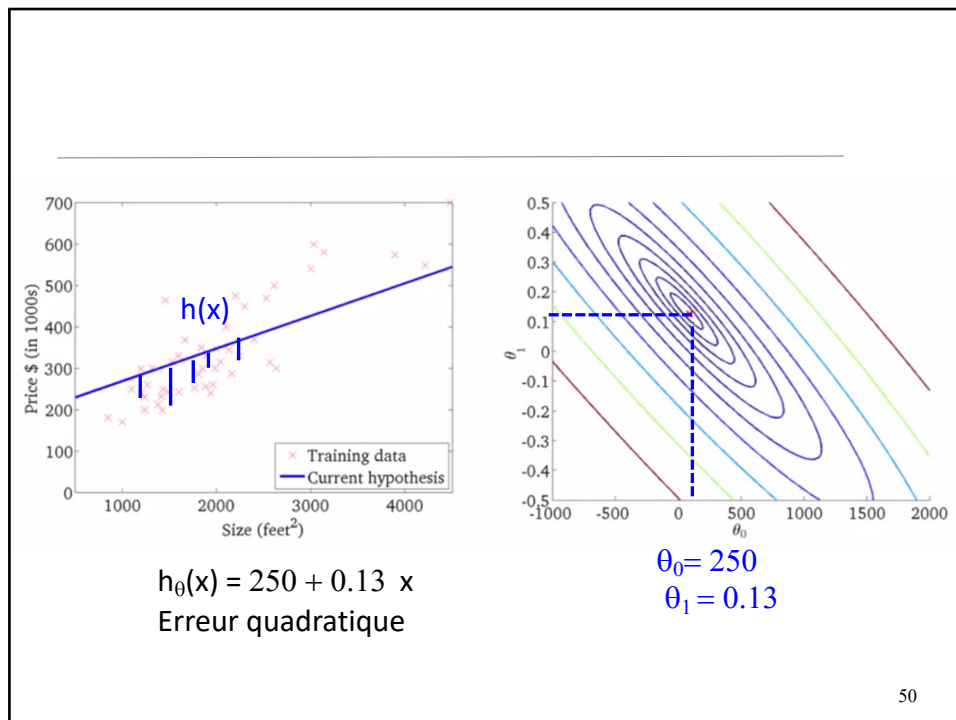
47



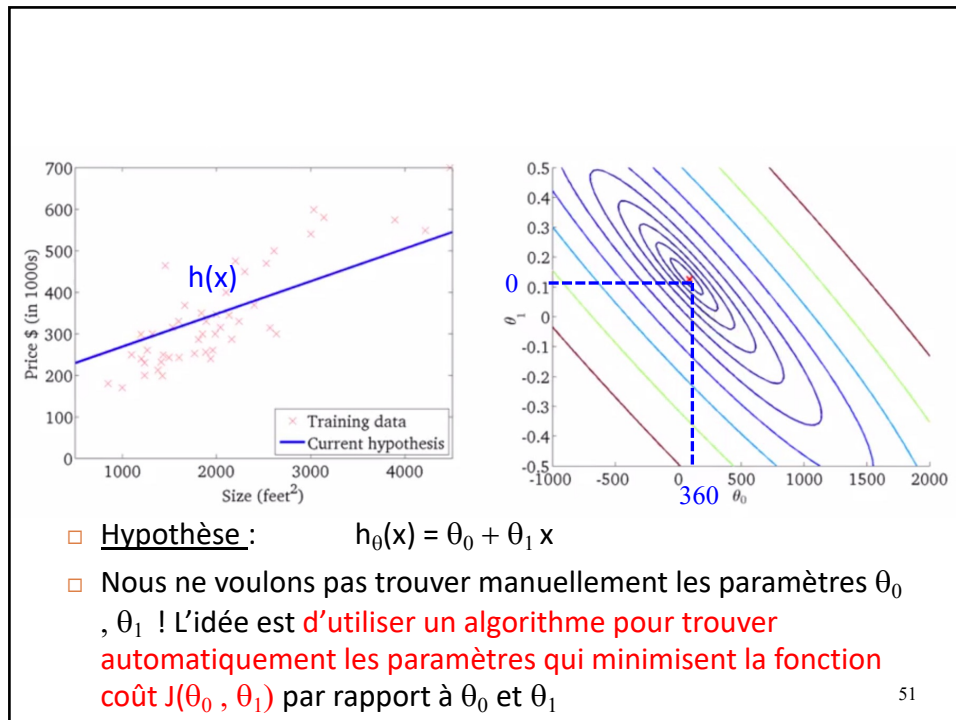
48



49



50



51

Université
Bretagne Sud
ubs:

Régression linéaire Descente de gradient

52

52

- Fonction coût : $J(\theta_0, \theta_1)$

□ But : minimiser $J(\theta_0, \theta_1)$
 θ_0, θ_1

$J(\theta_0, \theta_1, \theta_2, \dots \theta_n)$

minimiser $J(\theta_0, \theta_1, \theta_2, \dots \theta_n)$
 $\theta_0, \theta_1, \theta_2, \dots \theta_n$

53

53

- Fonction coût : $J(\theta_0, \theta_1)$

□ But : minimiser $J(\theta_0, \theta_1)$
 θ_0, θ_1

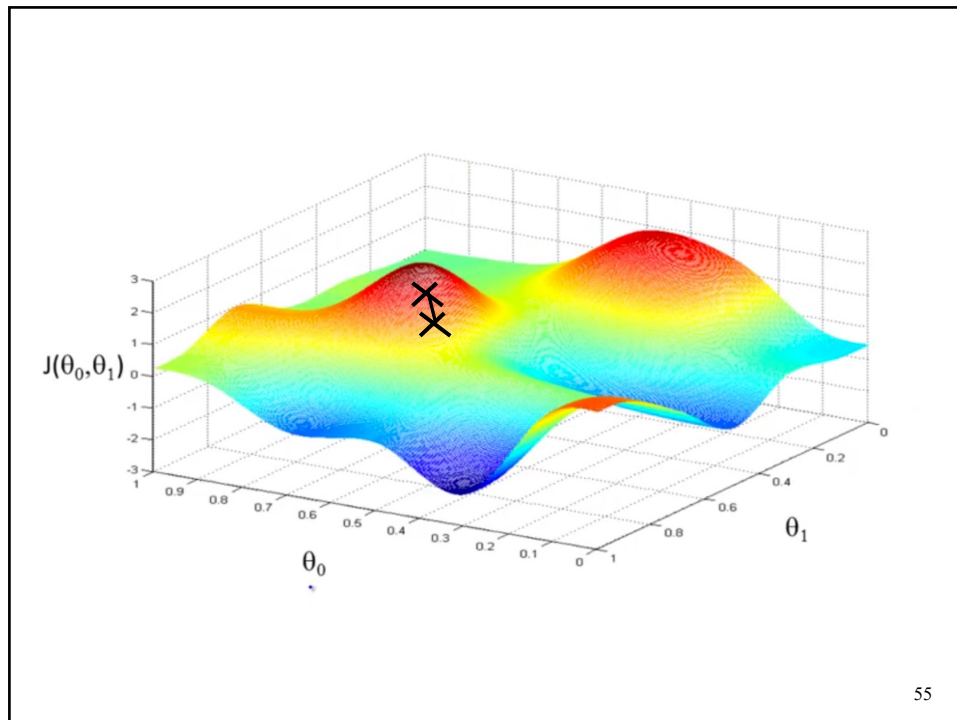
$J(\theta_0, \theta_1, \theta_2, \dots \theta_n)$

minimiser $J(\theta_0, \theta_1, \theta_2, \dots \theta_n)$
 $\theta_0, \theta_1, \theta_2, \dots \theta_n$
- Algorithme :

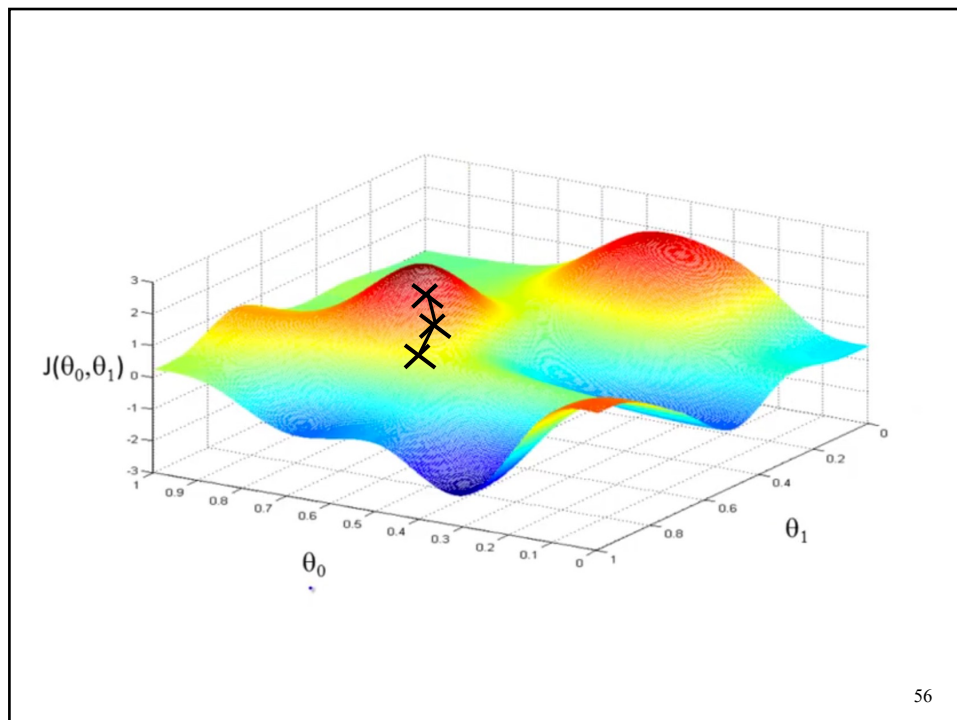
 - Commencer avec θ_0, θ_1 (initial guesses, e.g., $\theta_0 = 0, \theta_1 = 0$)
 - Changer θ_0, θ_1 pour réduire $J(\theta_0, \theta_1)$
jusqu'à atteindre un minimum

54

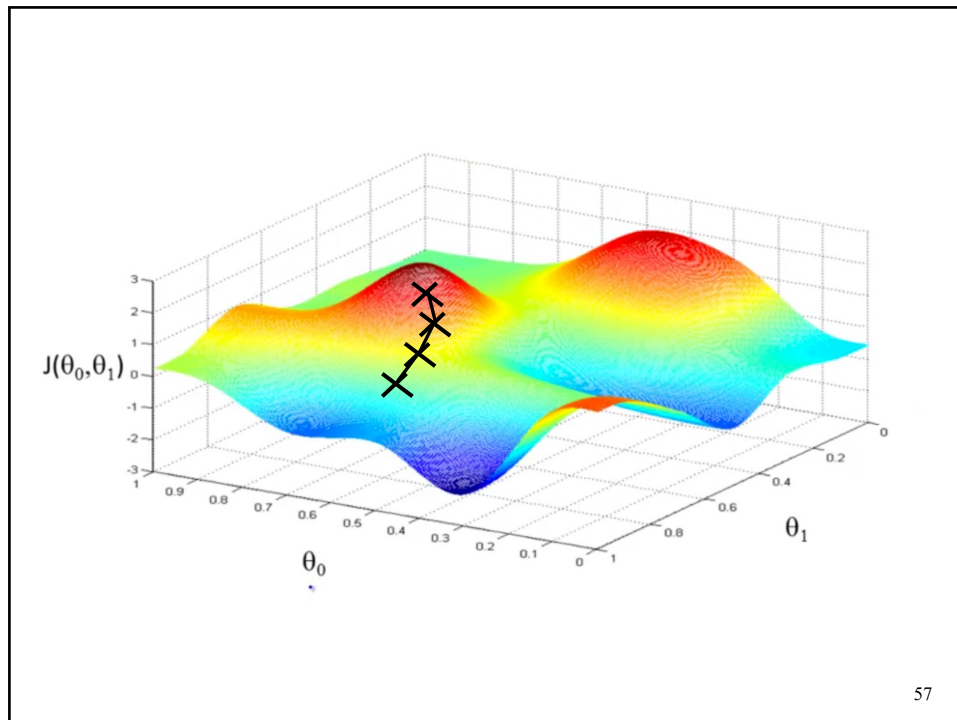
54



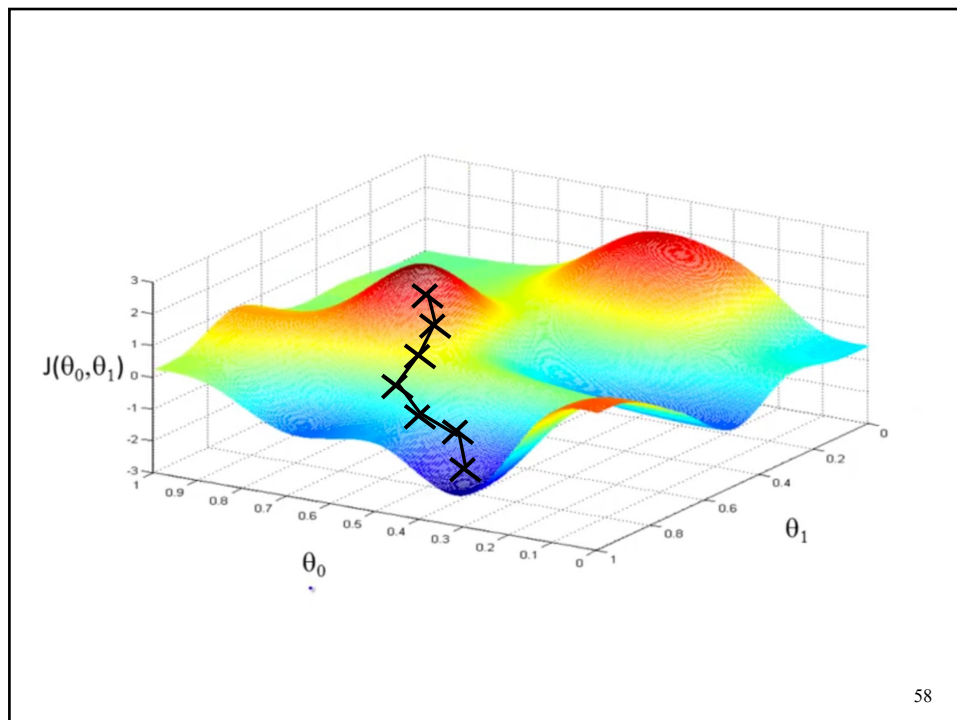
55



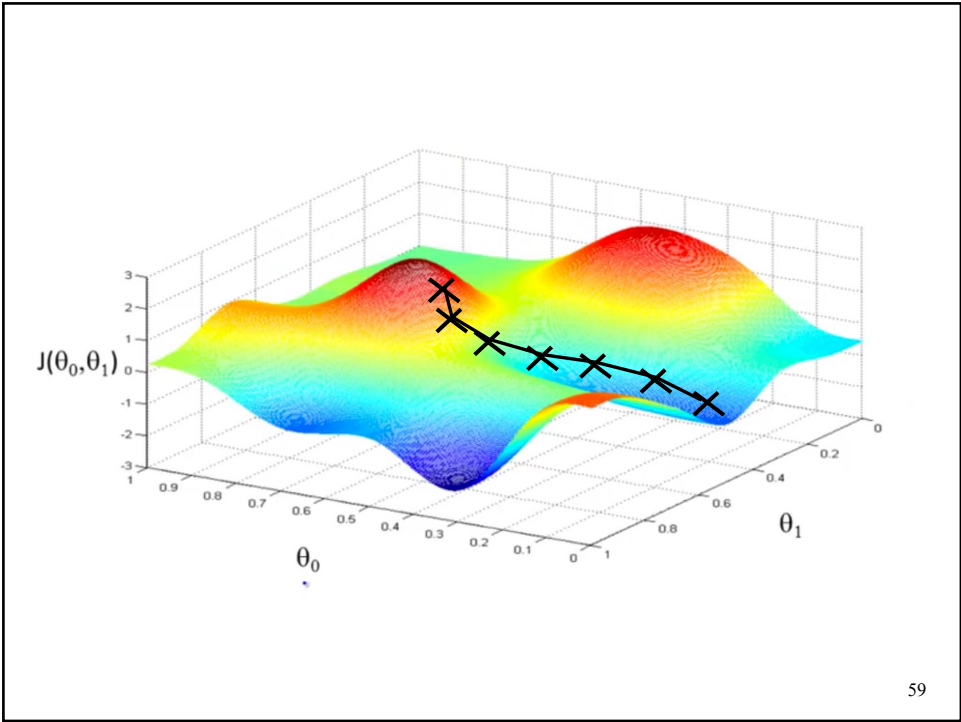
56



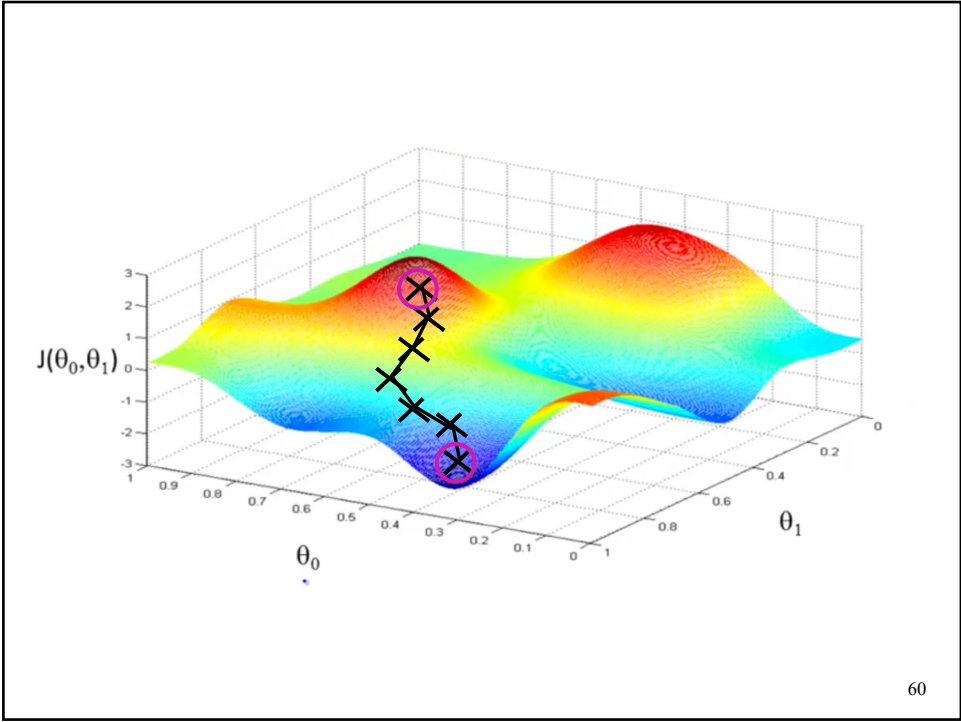
57



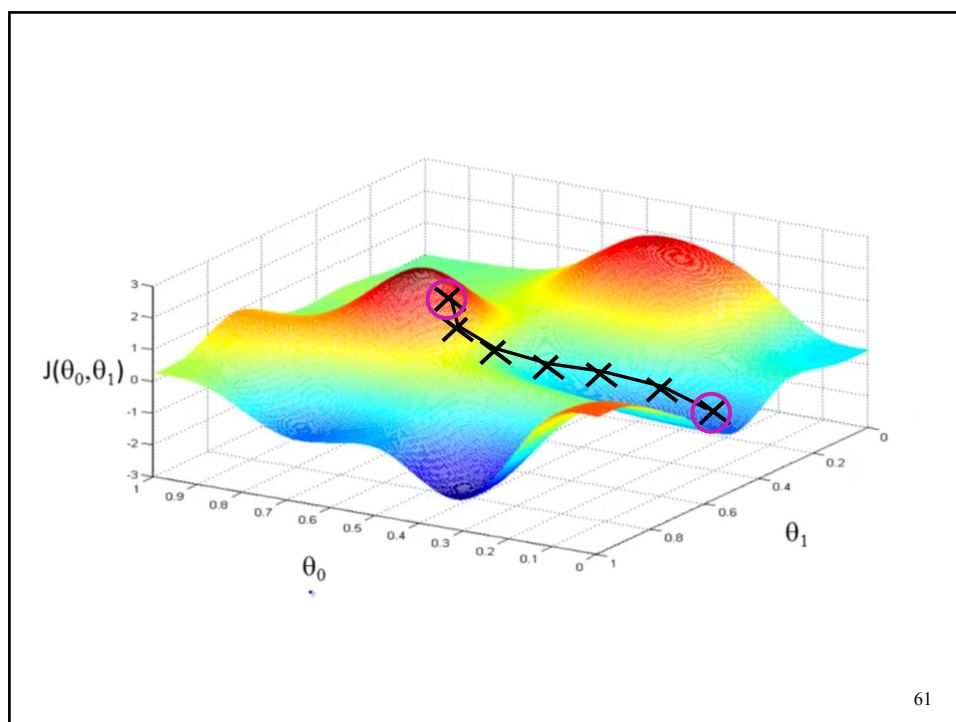
58



59



60



61

Algorithme de descente de gradient

▣ Algorithme :

Répéter jusqu'à convergence :

Repeat {

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

until θ_0, θ_1 converge

}

Mises à jour simultanées de θ_0, θ_1

62

62

Gradient Descent Algorithm

Algorithm:

Répéter jusqu'à convergence :

Repeat {

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

}

Paramètre α : taux d'apprentissage

Contrôle la façon dont on prend la descente de plus grande pente (plus ou moins vite)

Mises à jour simultanées

$\text{temp}_0 = \theta_0 - \alpha \cdot \text{dérivée}(\theta_0, \theta_1)$

$\text{temp}_1 = \theta_1 - \alpha \cdot \text{dérivée}(\theta_0, \theta_1)$

$\theta_0 = \text{temp}_0$

$\theta_1 = \text{temp}_1$

Implémentation fausse

$\text{temp}_0 = \theta_0 - \alpha \cdot \text{dérivée}(\theta_0, \theta_1)$

$\theta_0 = \text{temp}_0$

$\text{temp}_1 = \theta_1 - \alpha \cdot \text{dérivée}(\theta_0, \theta_1)$

$\theta_1 = \text{temp}_1$

63

63

Exercice : Supposez que $\theta_0=1$, $\theta_1=2$, et qu'on met à jour simultanément θ_0 et θ_1 en utilisant la règle : $\theta_i = \theta_i + \text{rac}(\theta_0 \theta_1)$.

Quelles sont les valeurs résultantes de θ_0 et θ_1 ?

- $\theta_0 = 1$ and $\theta_1=2$
- $\theta_0 = 1 + \text{rac}(2)$ and $\theta_1=2 + \text{rac}(2)$
- $\theta_0 = 2 + \text{rac}(2)$ and $\theta_1=1 + \text{rac}(2)$
- $\theta_0 = 1 + \text{rac}(2)$ and $\theta_1=2 + \text{rac}((1 + \text{rac}(2)).2)$

64

64

Descente du gradient Intuition

65

65

Algorithme de descente du gradient

▣ Algorithm:

Répéter jusqu'à convergence

Repeat {

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

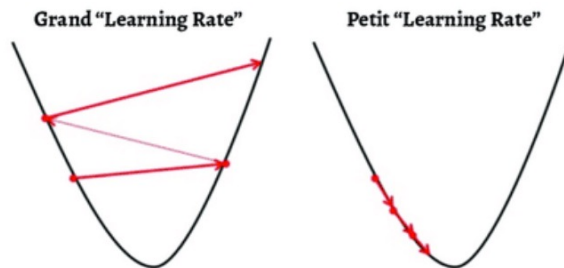
Mise à jour simultanée de θ_0 et θ_1

Supposons que nous ayons
uniquement θ_1

66

66

Algorithme de descente du gradient



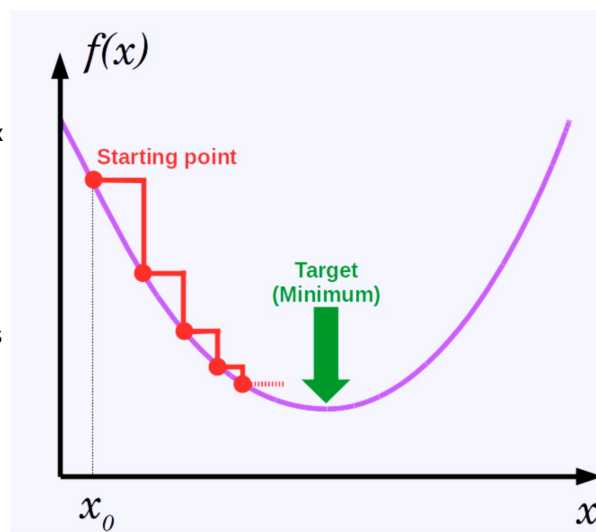
- Si α est trop grand, la descente de gradient peut "overshooter" le minimum (osciller de part et d'autre). l'algorithme peut ne pas converger, ou même diverger.
- Si α est trop petit, la descente de gradient peut être (trop) lente

67

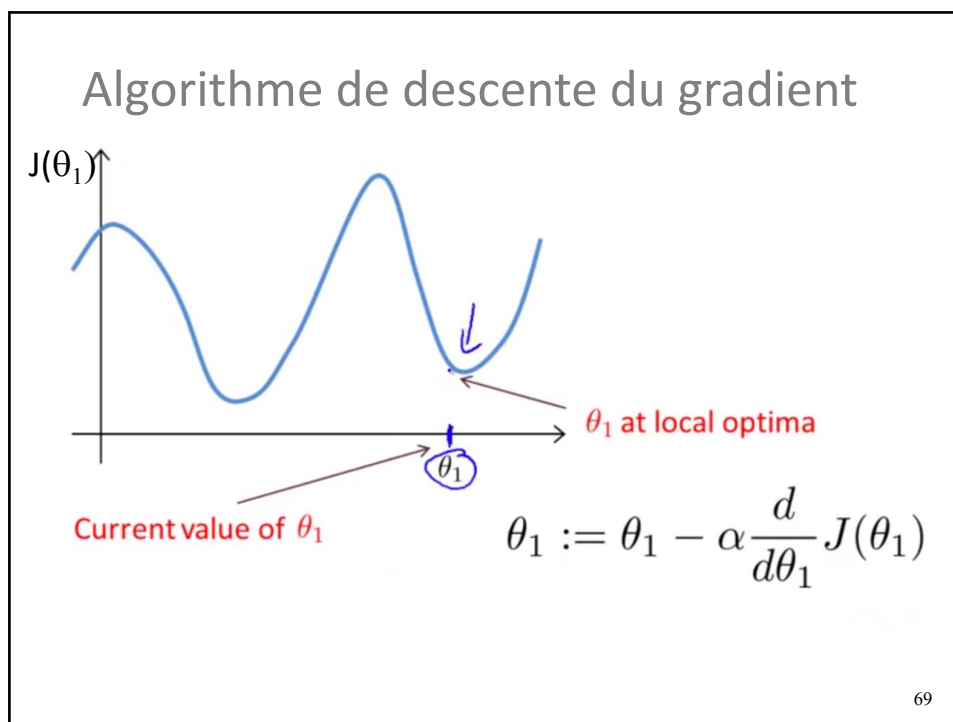
67

Algorithme de descente du gradient

- La descente de gradient peut converger vers un minimum local, avec un taux α fixé.
- Lorsqu'on approche le minimum local, la descente de gradient va prendre automatiquement de plus petits incréments : ainsi, pas besoin de diminuer α au cours du temps.



68

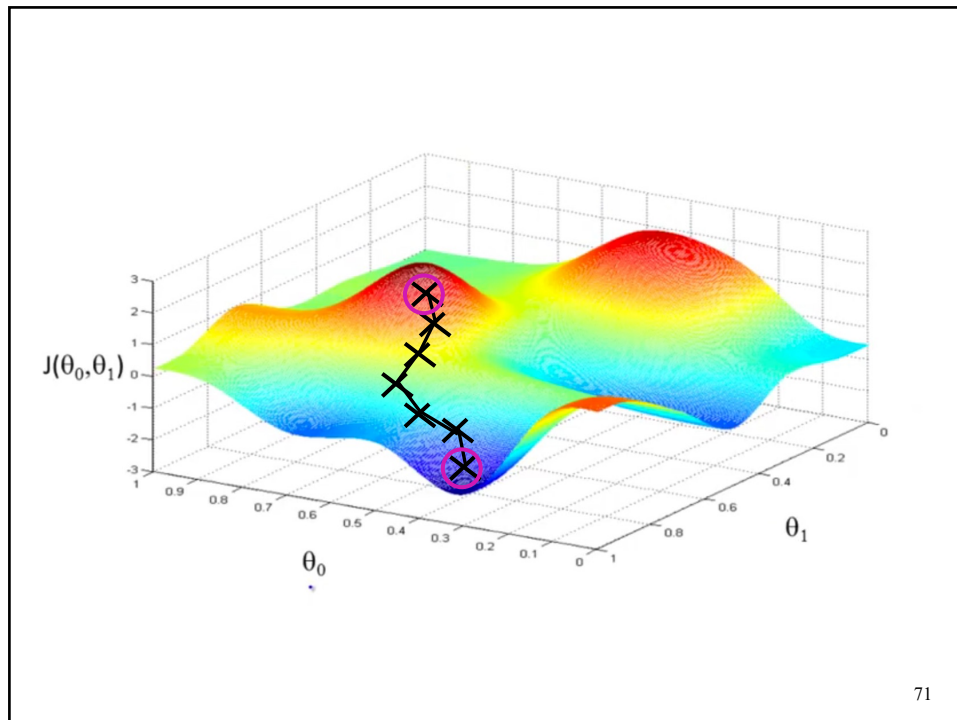


69

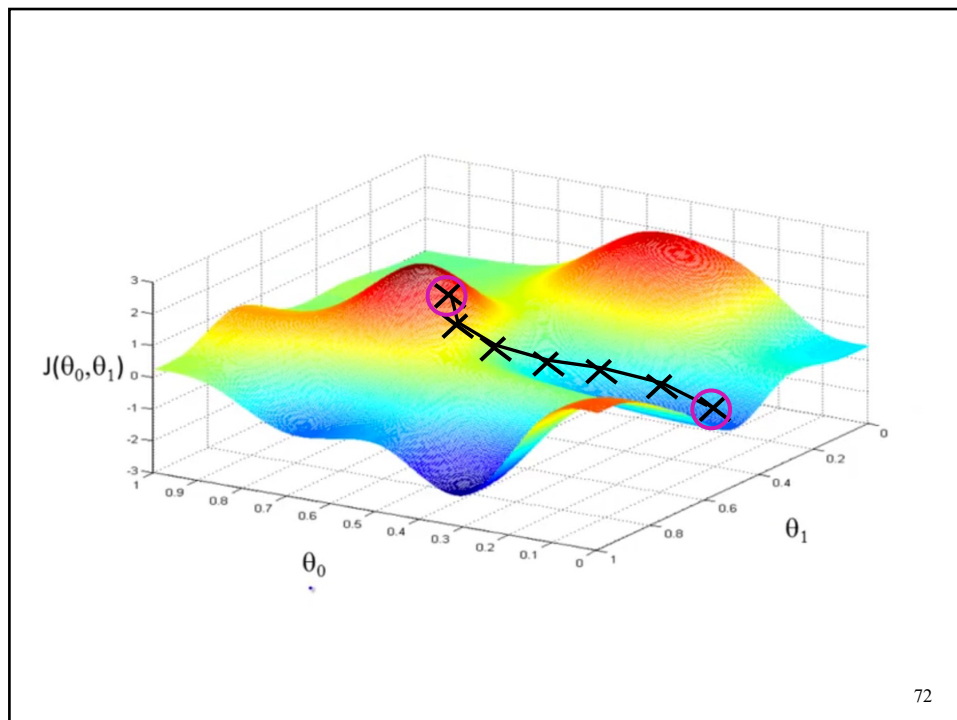
Descente du gradient pour la régression linéaire

70

70

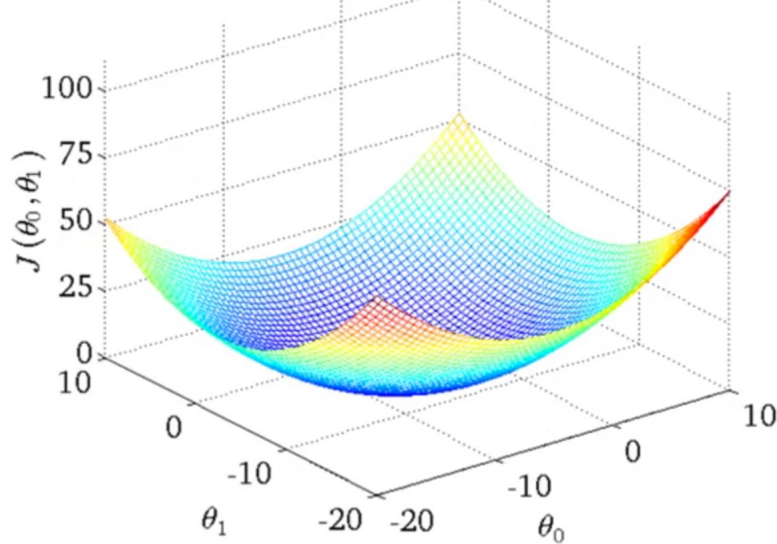


71



72

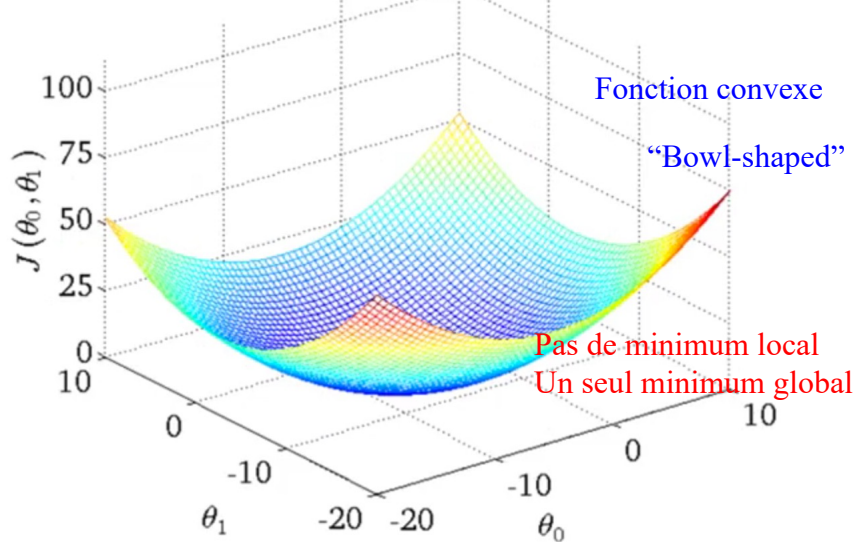
Régression : coût quadratique



73

73

Régression : coût quadratique



74

74

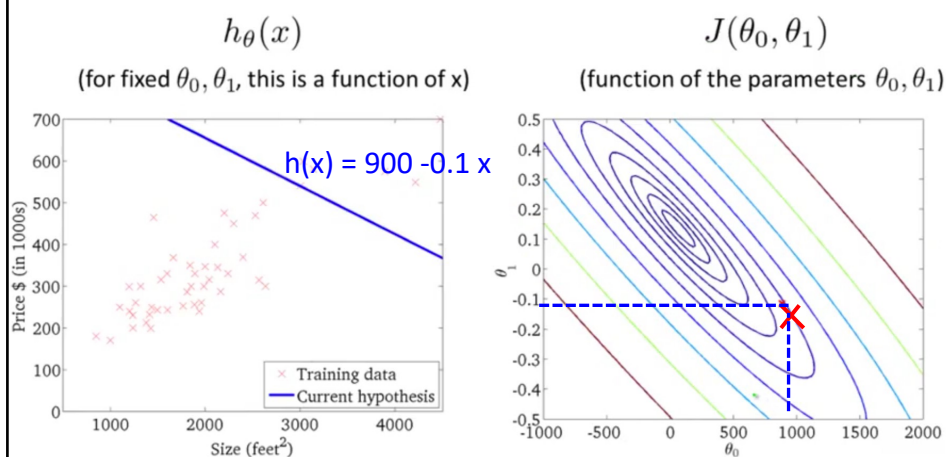
Régression : coût quadratique

- **Algorithme de descente de gradient** : peut tomber dans un minimum local, dépendant du point de départ de l'algorithme (initialisation des paramètres θ_0 et θ_1).
- **Régression linéaire** : la descente de gradient appliquée à une fonction coût quadratique :
 - La fonction coût est toujours une fonction convexe (bowl-shaped)
 - Avec une telle méthode, impossible de tomber dans un minimum local, il n'y a qu'un seul minimum global.

75

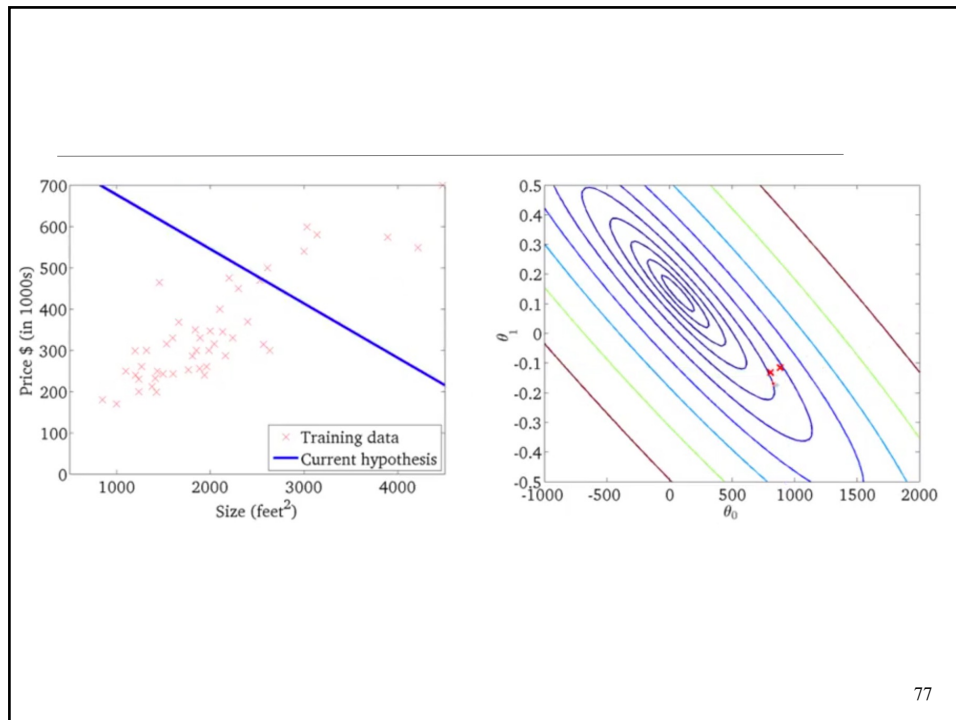
75

Algorithme en action

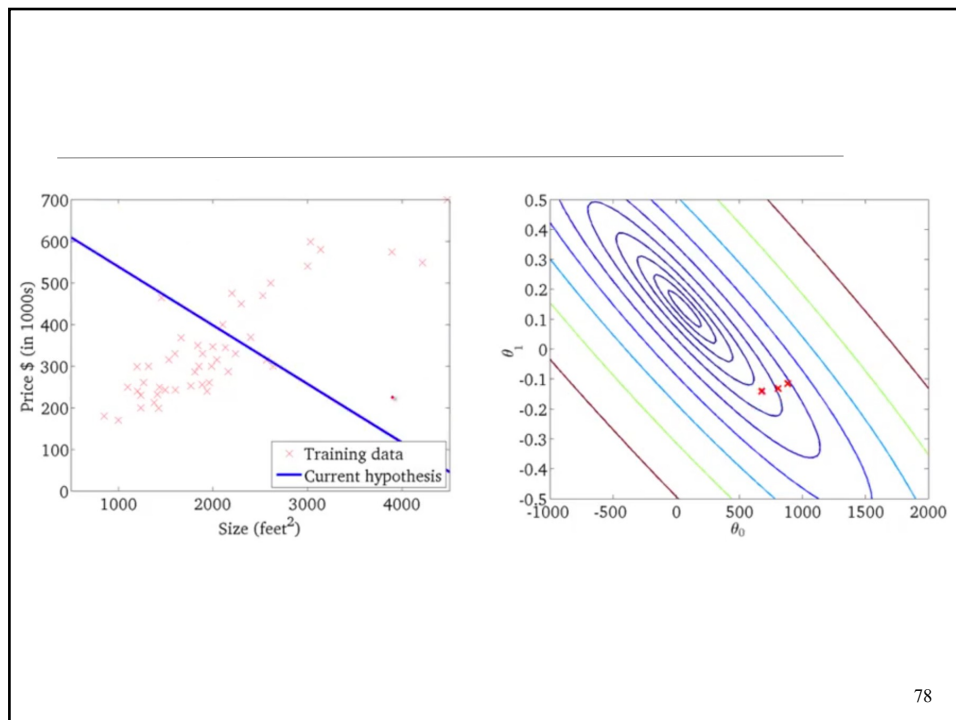


76

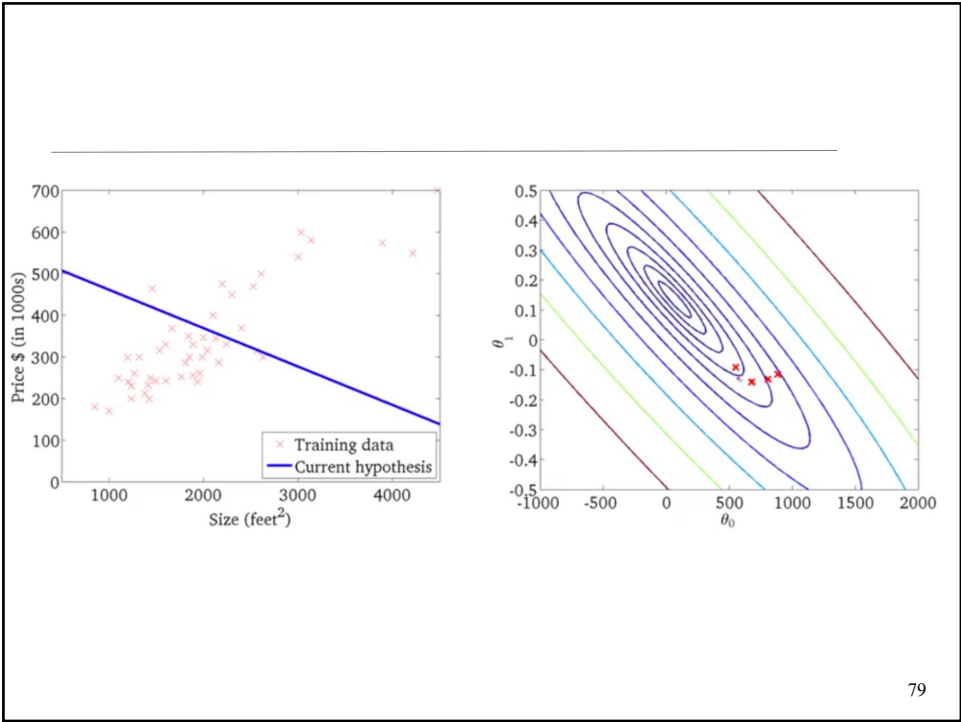
76



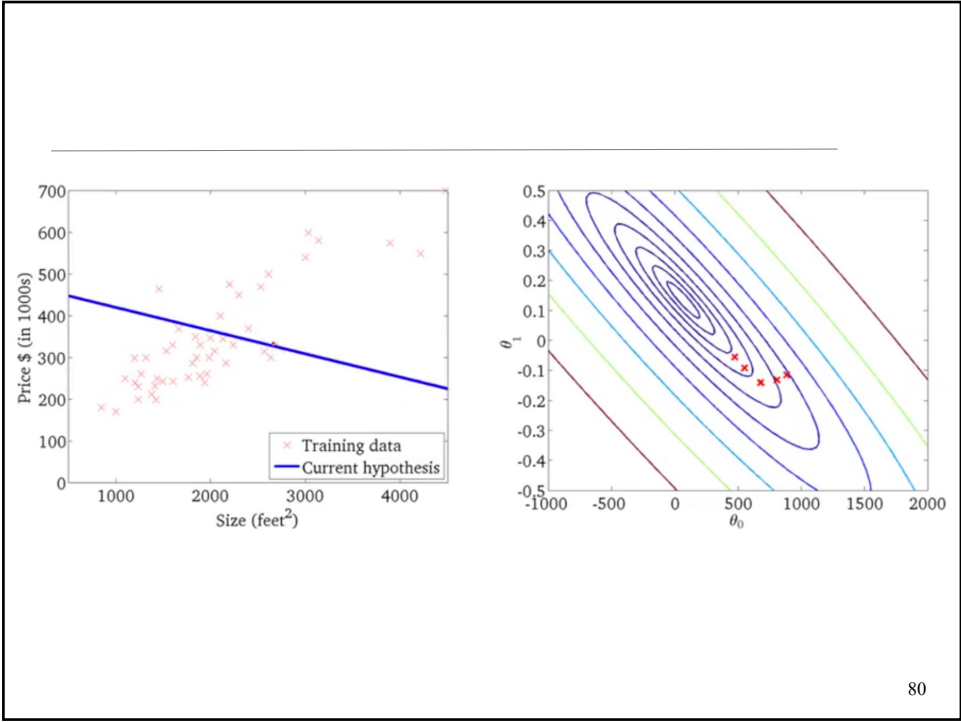
77



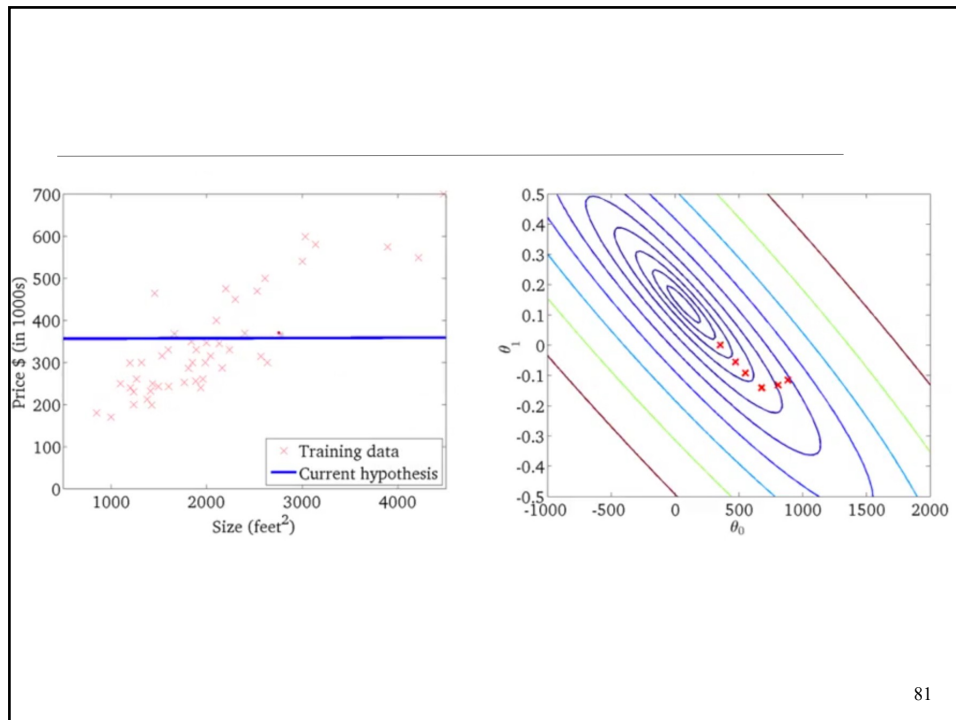
78



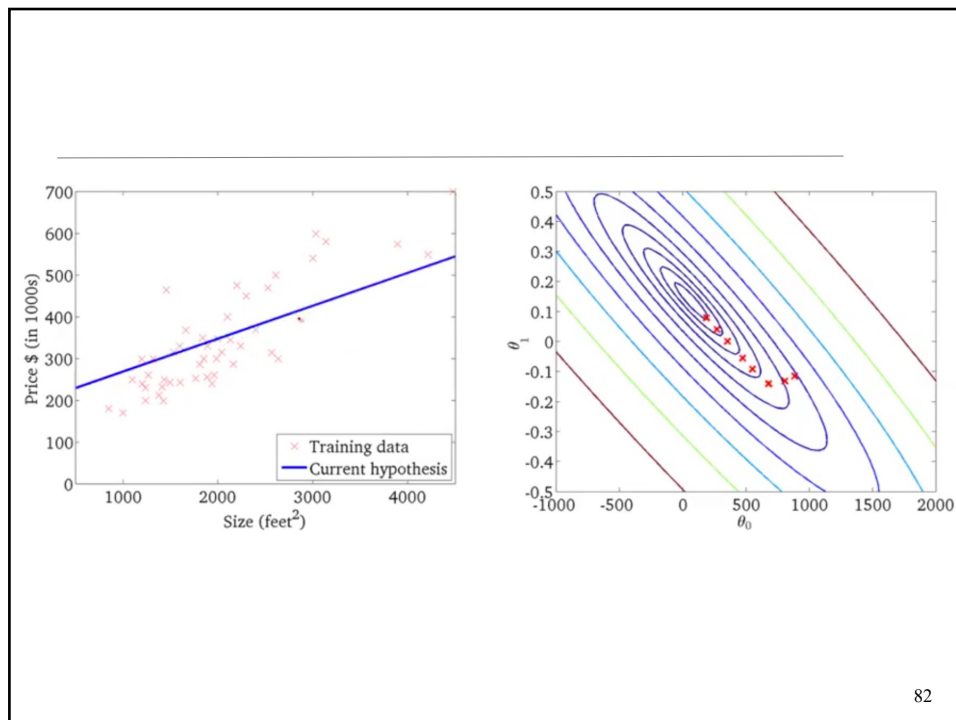
79



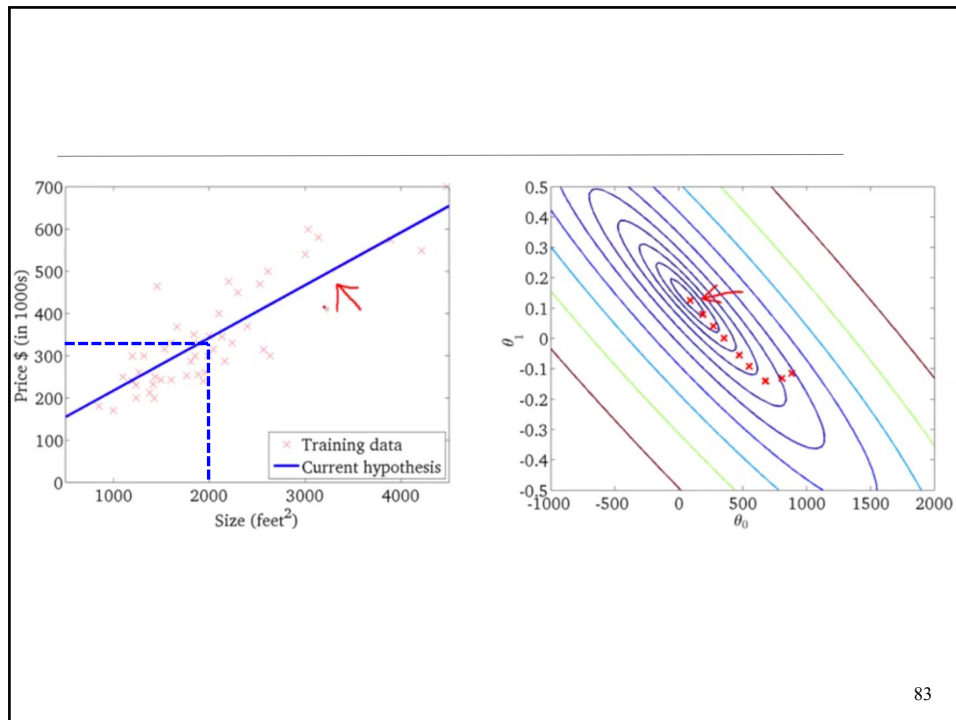
80



81



82



83

Parmi les propositions suivantes, lesquelles sont vraie ?
Sélectionnez toutes les bonnes réponses

- ☐ Pour que la descente de gradient converge, il faut que le taux α décroisse au cours du temps.
- ☐ La descente de gradient trouve toujours le minimum global de toute fonction coût $J(\theta_0, \theta_1)$.
- ☐ La descente de gradient peut converger même si α est maintenu fixe (mais α ne peut pas être trop grand, sinon l'algorithme peut ne pas converger).
- ☐ Pour le choix spécifique de fonction coût utilisé dans la régression linéaire, il n'y a pas d'optimum local (autre que l'optimum global).

84

<p>▣ <u>Algorithm de descente du gradient</u></p> <hr/> <p>Repeat {</p> $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$ <p>(for j = 1 and j = 0) until convergence }</p>	<p>▣ <u>Linear Regression Model</u></p> <hr/> $h_\theta(x) = \theta_0 + \theta_1 x$ $J(\theta_0, \theta_1) = \frac{1}{2m} \cdot \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$ $= \frac{1}{2m} \cdot \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$ <p>▣</p>
--	--

85

<p>Algorithme de descente du gradient</p> <hr/> $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \cdot \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$ $= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \cdot \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$ $j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \cdot \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$ $j = 1 : \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \cdot \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$
--

86

Algorithme de descente du gradient

Repeat until convergence {

$$j = 0 : \theta_0 = \theta_0 - \frac{1}{m} \cdot \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

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

} ➔ update θ_0 and θ_1 simultaneously

87

87