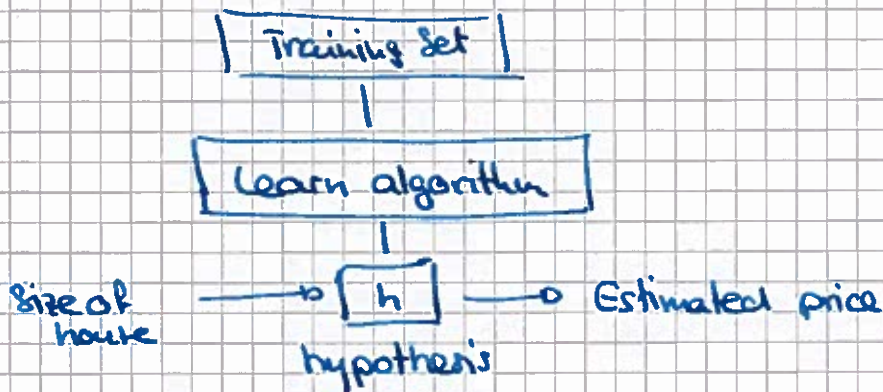


Machine learning  
- Stanford University

## 2. Linear Regression

## Notation:

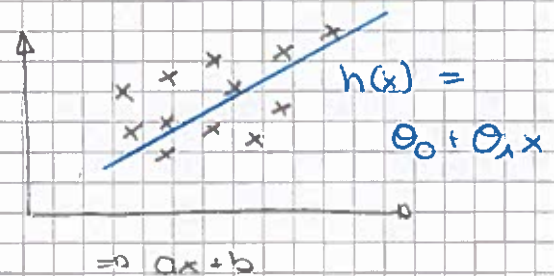
- $m$  - number of training examples
- $x$  - input variable / feature
- $y$  - output variable
- $(x, y)$  simple training example
- $(x^{(i)}, y^{(i)})$  specific example



## Representation of $h$

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

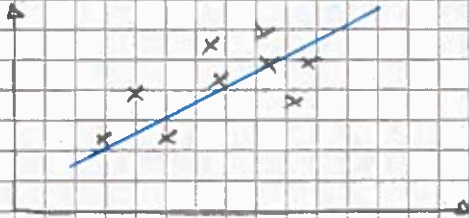
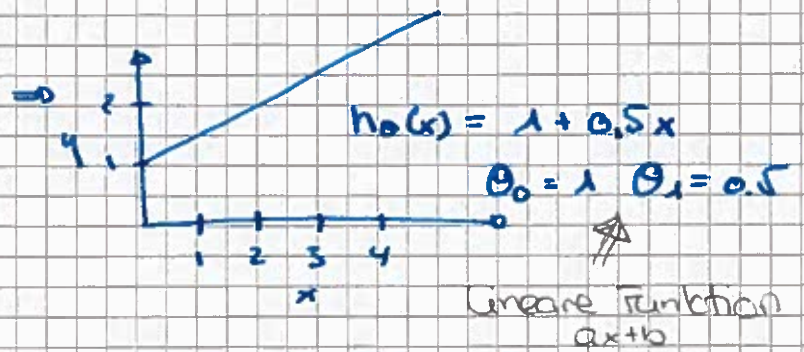
↳ shorthand  $h(x)$



Linear Regression with one variable  
Univariate Linear Regression



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



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

$$\text{minimize}_{\theta_0, \theta_1} \frac{1}{2m} \sum_{i=1}^m \underbrace{(h_{\theta}(x^{(i)}))}_{\text{Predicted value}} - \underbrace{y^{(i)}}_{\text{real value}})^2$$

# of training examples

||

$$h(x^{(i)}) = \theta_0 + \theta_1 x \Rightarrow \text{given by the data}$$

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

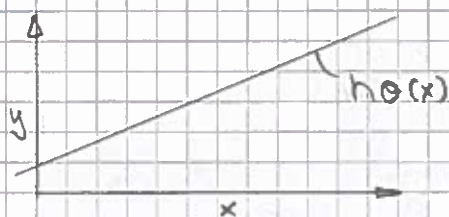
$$\text{minimize} = J(\theta_0, \theta_1) \Rightarrow \text{Cost Function}$$

$$\Rightarrow \text{Squared Error Function}$$

Summary:

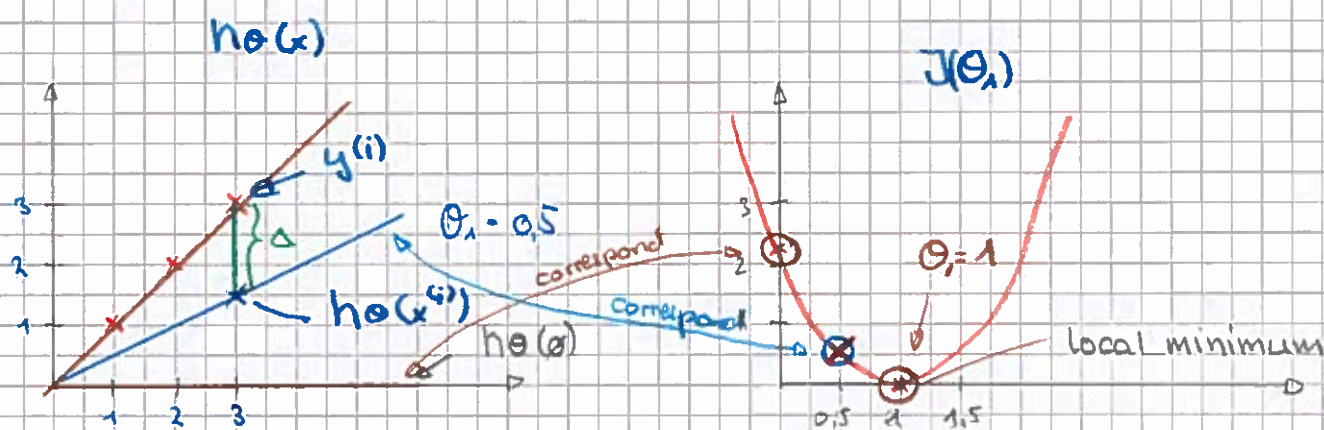
Hypothesis:  $h(x) = \theta_0 + \theta_1 x$

Parameters:  $\theta_0, \theta_1$



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

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



$$J(0.5) = \frac{1}{2m} \sum_{i=1}^m = \frac{1}{2 \cdot 3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2]$$

$$= \frac{1}{2 \cdot 3} \cdot 3.5 = 0.58$$

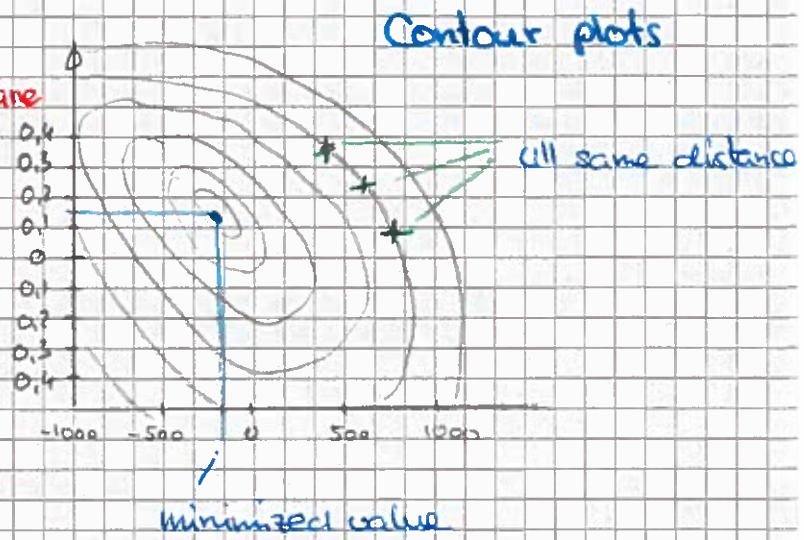
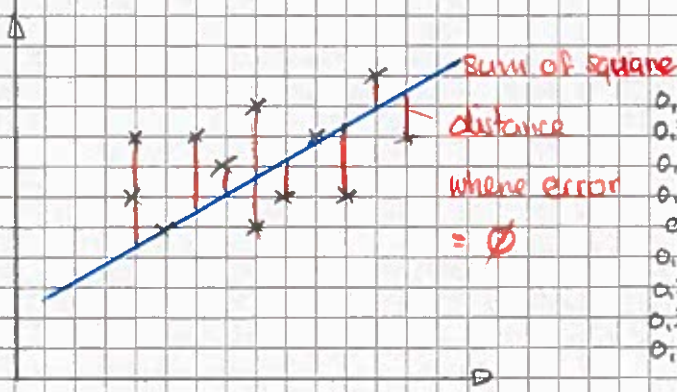
$\Delta$  difference between predicted and effective value

minimize  $J(\theta_1)$

$$J(0) = \frac{1}{2m} \cdot [(0 - 1)^2 + (0 - 2)^2 + (0 - 3)^2] = \frac{1 \cdot 14}{2 \cdot 3} = \sim 2.3$$

$$J(1) = \frac{1}{2m} \cdot [(1 - 1)^2 + (2 - 2)^2 + (3 - 3)^2] = \frac{0}{6} = 0$$





# Gradient descent algorithm

5

=> want to get the local/global minimum

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

=> Gradient descent applies to more general functions

$$J(\theta_0, \theta_1, \dots, \theta_n) \Rightarrow \min J(\theta_0, \theta_1, \dots, \theta_n)$$

Repeat until convergence  $\epsilon$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j=0, j=1)$$

=> simultaneous update

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

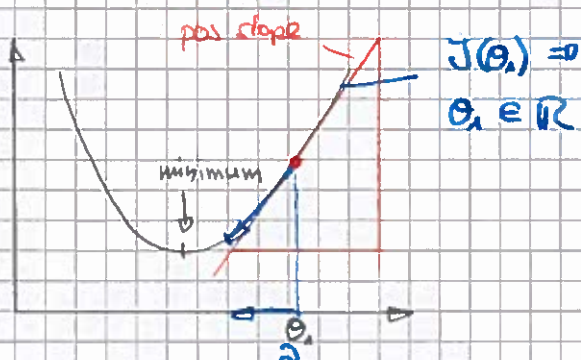
$\alpha$  = learning rate

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

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

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

derivative term

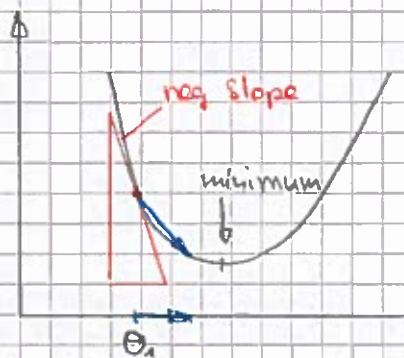


$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

$\geq 0$

$$\theta_1 := \theta_1 - \alpha \cdot (\text{pos. number})$$

always pos.



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

$\leq 0$

$$\theta_1 := \theta_1 - \alpha \cdot (\text{neg. number})$$

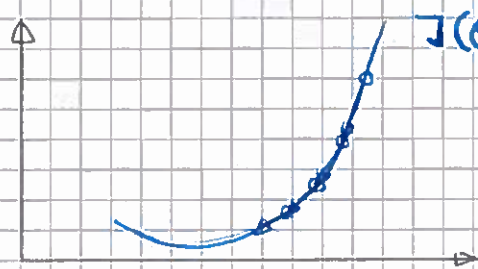
make sure that  $\alpha$  is not too small or too large

- too small might take too many steps to converge
- too big might overshoot the min. and might never converge.





Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed.



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

As we approach local minimum, gradient descent will take smaller steps.

term  $x$  will decrease therefore value will be smaller & smaller

No need to decrease  $\alpha$ .

## Linear Regression with One Variable

Apply gradient descent to linear regression function

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

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

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

## Gradient descent algorithm

repeat until convergence  $\Sigma$

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

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

}

always update simultaneously.