# 2.   LINEAR REGRESSION
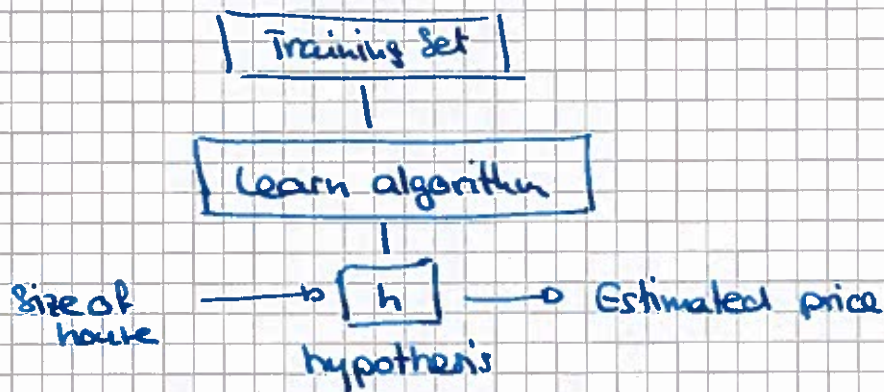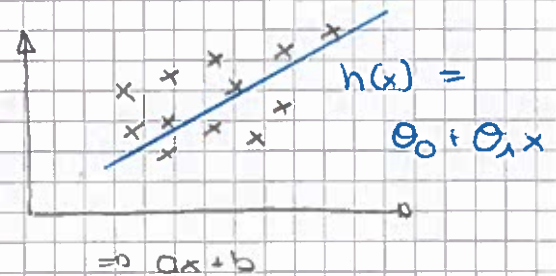
# Chapter 2

# Linear Regression

Notation:

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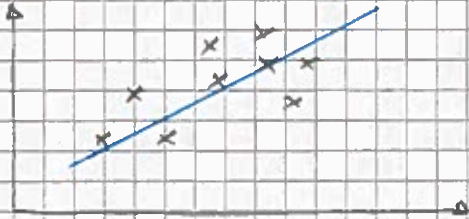$(x^{(i)}, y^{(i)})$  specific example

| Training Set |

| Learn algorithm |

Size of house $\longrightarrow$ | h | $\longrightarrow$ Estimated price
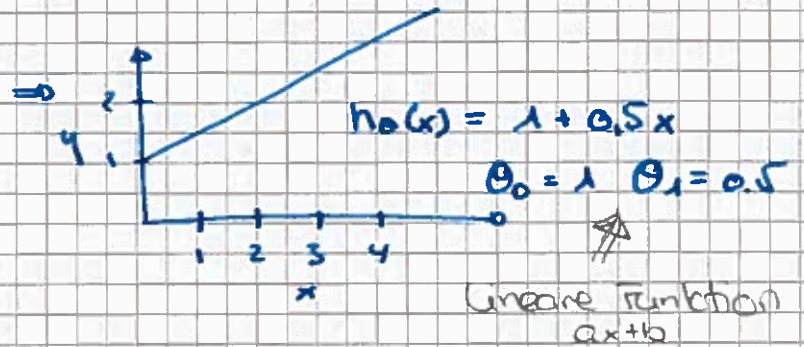
hypothesis

Representation of h

$\Rightarrow$  $h\theta(x) = \theta_0 + \theta_1 x$  $\Rightarrow$

$\rightarrow$ shorthand $h(x)$

Linear Regression with one variable
Univariate Linear Regression

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

$\Rightarrow ax+b$

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

$\Rightarrow$

$$h_\theta(x) = 1 + 0.5x$$
$$\theta_0 = 1 \quad \theta_1 = 0.5$$

Lineare Funktion
$ax + b$

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

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

\# of training examples

Predicted value

real value

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

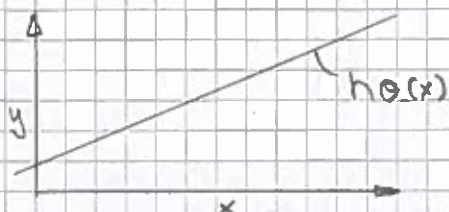$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

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

$$\Rightarrow \text{Squared error function}$$

Summary:

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

Parameters :    $\theta_0, \theta_1$

Cost function :    $J(\theta_0, \theta_1) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} \left( h(x^{(i)}) - y^{(i)} \right)^2$

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

$J(0,5) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} \overset{=0}{} \left[ (0,5-1)^2 + (1-2)^2 + (1,5-3)^2 \right]$

$\overset{=0}{} \dfrac{1}{2 \cdot 3} \cdot 3.5 = 0,58$

$\triangle$  difference between predicted and effective value

minimize  $J(\theta_1)$

$J(0) \Rightarrow \dfrac{1}{2m} \cdot \left[ (0-1)^2 + (0-2)^2 + (0-3)^2 \right] = \dfrac{1 \cdot 14}{2 \cdot 3} = \sim 2,3$

$J(1) \overset{=0}{} \dfrac{1}{2m} \cdot \left[ (1-1)^2 + (2-2)^2 + (3-3)^2 \right] = \dfrac{0}{6} = \emptyset$

Contour plots



Sum of square distance where error = Ø

cill same distance

minimized value

=o want to get the local/global minimum

$$\min \; J(\Theta_0, \Theta_1)$$

=o Gradient descent applies to more general functions

$$J(\Theta_0, \Theta_1 \ldots \Theta_n) \quad =\!\!o \;\; \min J(\Theta_0, \Theta_1 \ldots \Theta_n)$$

Repeat until convergence ε

$$\Theta_j := \Theta_j - \alpha \frac{\partial}{\partial_j} J(\Theta_0, \Theta_1) \quad (\text{for } j=0, \; j=1)$$

}

-o simultanious update

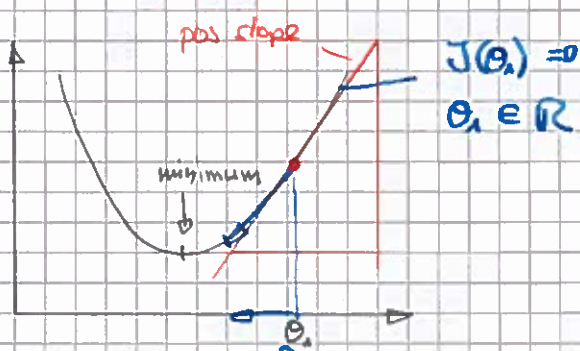$$temp\emptyset := \Theta_0 - \alpha \frac{\partial}{\partial \Theta_0} J(\Theta_0, \Theta_1)$$

$$temp1 := \Theta_1 - \alpha \frac{\partial}{\partial \Theta_1} J(\Theta_0, \Theta_1)$$

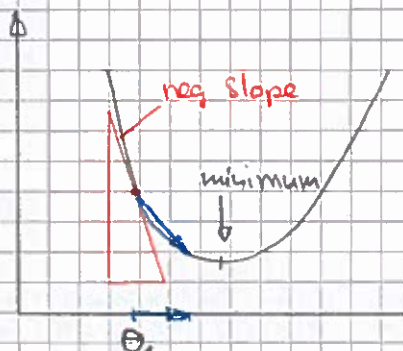$$\Theta_0 := temp\emptyset$$

$$\Theta_1 := temp1$$

$$\underbrace{\qquad\qquad}_{\text{derivative term}}$$

$\alpha$ = learning rate



$$J(\Theta_1) = 0$$
$$\Theta_1 \in \mathbb{R}$$

$$\Theta_1 := \Theta_1 - \alpha \underline{\frac{\partial}{\partial \Theta_1} \cdot J(\Theta_1)}$$

$$\geq 0$$

$$\underline{\frac{\partial}{\partial \Theta_1} \cdot J(\Theta_1)}$$

$$\leq \emptyset$$

$$\Theta_1 := \Theta_1 - \alpha \cdot (\text{pos. number})$$

always pos.

$$\Theta_1 := \Theta_1 - \alpha \cdot (\text{neg. number})$$
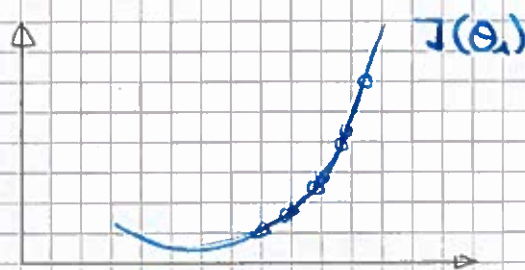
make sure that $\alpha$ is not too small or to large
- too small might take to many steps to converge
- too big might overshoot the min. and might never converge.

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



$$\theta_1 := \theta_1 - \alpha \underbrace{\frac{\partial}{\partial \theta_1} \cdot J(\theta_1)}_{\times}$$

term $\times$ will decrease therefore value will be smaller & smaller

As we approach local minimum, gradient descent will take smaller steps. No need to decrease $\alpha$.

### Linear Regression with One Variable

Apply gradient descent to linear regression function

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

$$= \frac{\partial}{\partial \theta_1} \cdot \frac{1}{2m} \sum_{i=1}^{m} \left( \theta_0 + \theta_1 x^{(i)} - y^{(i)} \right)^2$$

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

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

### Gradient descent algorithm

repeat until convergence {

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

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

}

always update simultaniously.