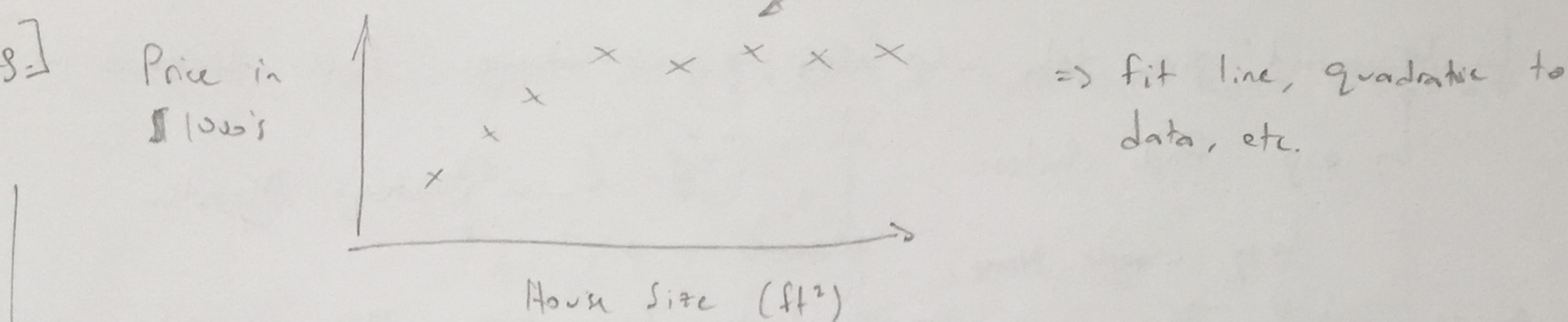


## COURSERA ML - WEEK 1

### What Is Machine Learning?

- Supervised Learning:

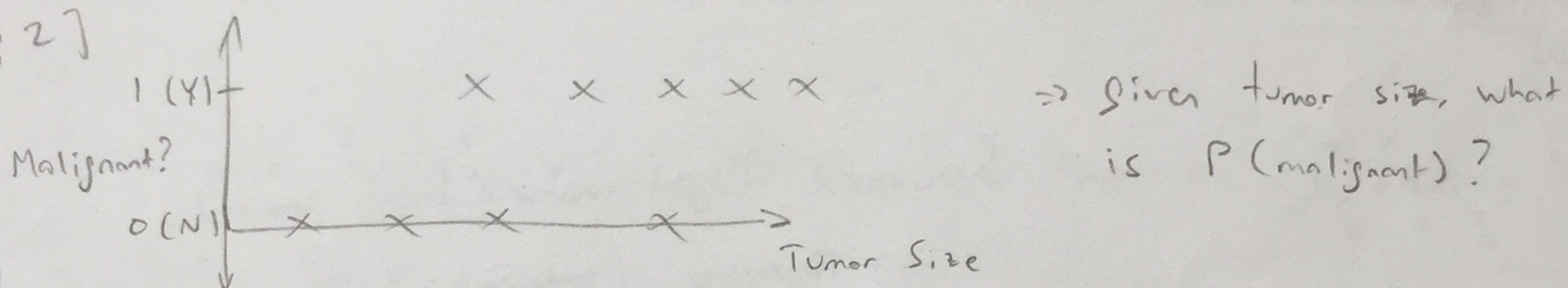
e.g.] Price in  
\$1000's



↳ SUPERVISED → RIGHT ANSWERS GIVEN (fit algorithm to data set)

↳ Regression Problem: predict continuous valued output (e.g. price)

e.g 2]



↳ CLASSIFICATION PROBLEM → predict discrete output

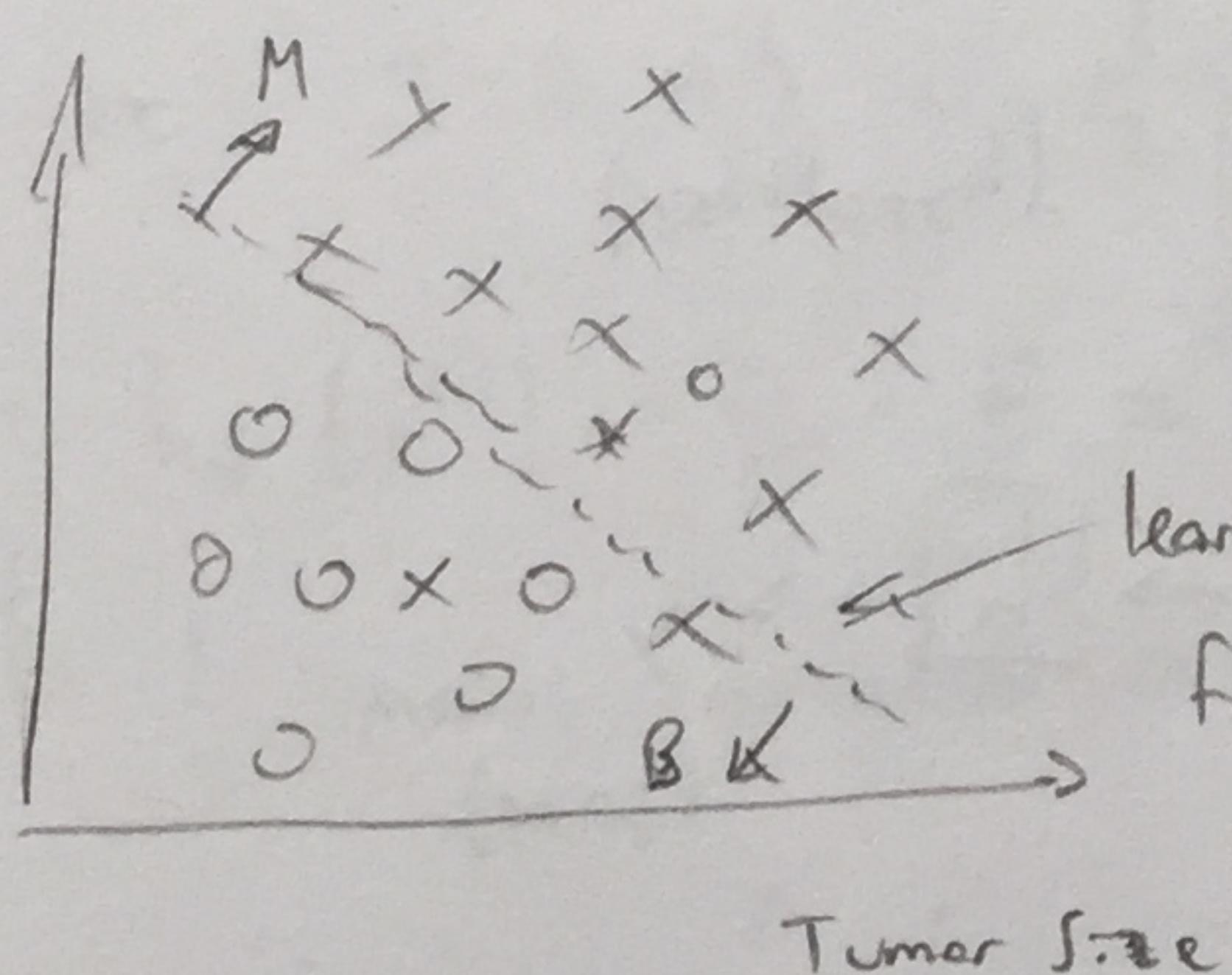
↳ can have 2, 3, ... discrete outcomes

REPLOT:

~~— x — x — o — x — o — x — o — >~~

e.g 3] Multiple Features / Attributes:

Age



learning algorithm

fits line

Some algorithms need ∞ features

(math trick)

- Regression → CONTINUOUS VALUE OUTPUTS

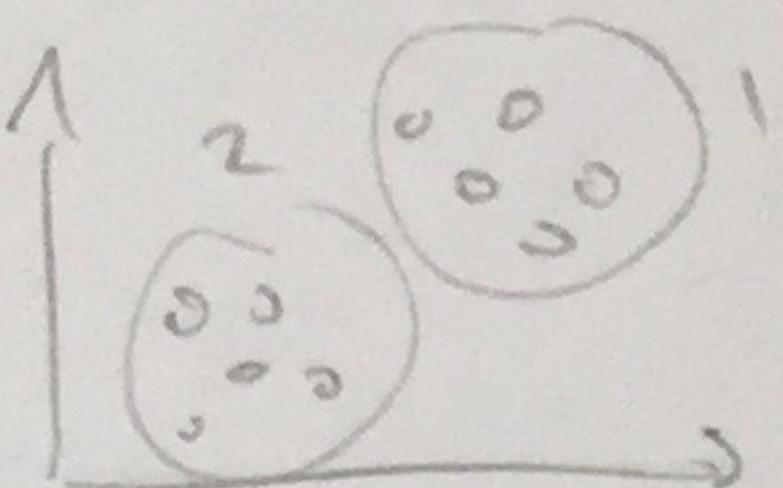
- Classification → DISCRETE OUTPUTS

## • Unsupervised Learning

↳ Supervised: given data set, get discrete outputs (classification) or  
continuous outputs (regression)

↳ Unsupervised: given data w/o labels → make inferences / find structure in data

e.g. find clusters



→ Clustering Algorithm

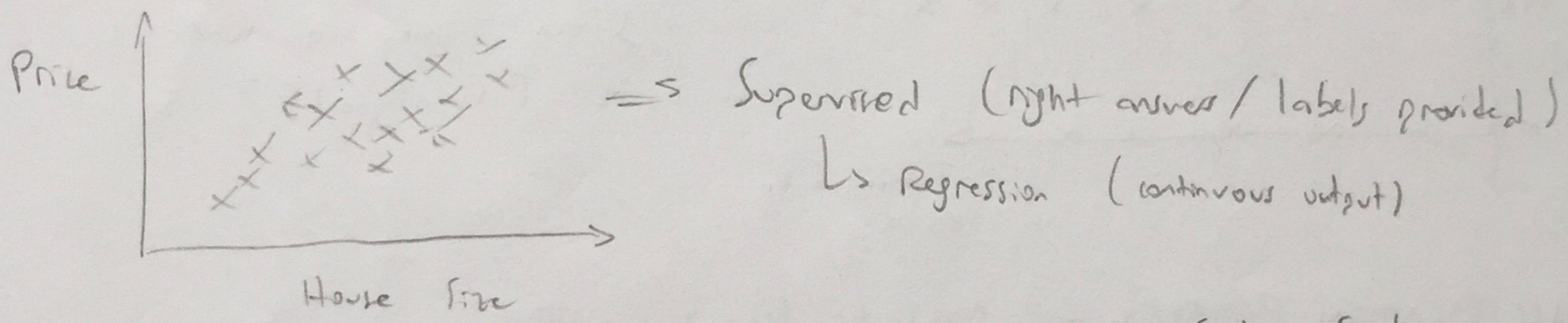
↳ c. Google News,

Genomics, Computer Clusters, Social Media, Market Segmentation, Astronomical

↳ e.g. cocktail problem → separate voices talking over each other

⇒ We are using Octave.

## Model Representation:



↳ Training set!  $x^{(i)}$  →

Size	Price
:	:
:	:
:	:
$i$	
:	:
:	:
:	:

Set of training examples  $(x, y)$

↓ # Rows → M

$x$ : input variable       $y$ : output variable

$(x^{(i)}, y^{(i)})$

↳  $i$ th training example

$\Rightarrow$  Training Set  $\rightarrow$  Learning Algorithm  $\rightarrow$   $\boxed{h}$  (hypothesis)

↳  $h$  maps  $x \rightarrow y \Leftrightarrow x \rightarrow [h] \rightarrow y$

↳ hypothesis written as  $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \dots$

Ls shorthand  $h(x)$

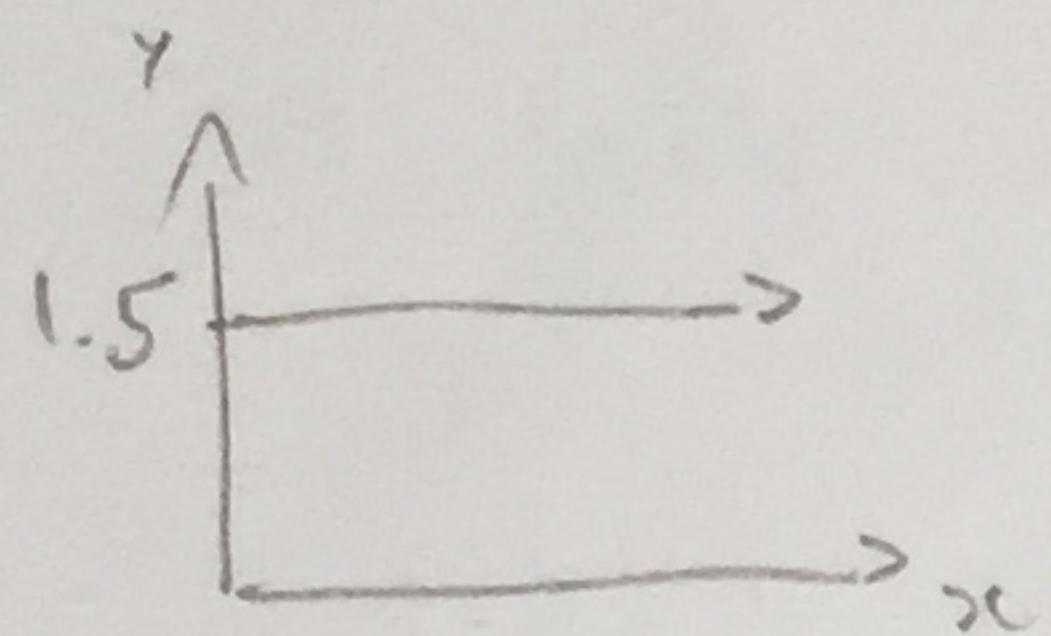
⇒ Linear Function:  $h_{\theta}(x) = \theta_0 + \theta_1 x$  ← linear regression w/ one variable ( $x$ )

↳ univariate (one variable) lin. reg.

### • Cost Function:

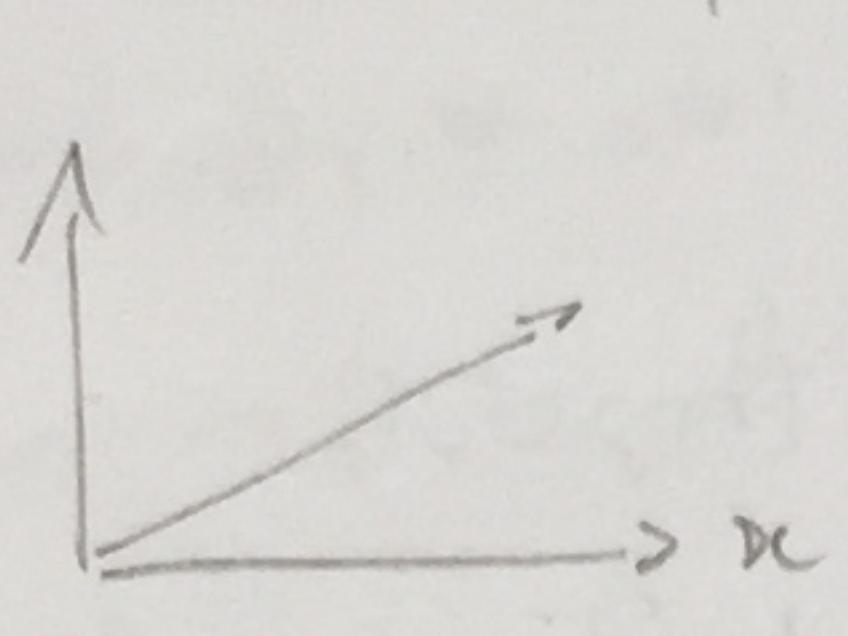
↳ We have training set  $(x, y) \times m$  & hypothesis  $h_{\theta}(x) = \theta_0 + x\theta_1$ .

↳ figure out  $\theta_0$  &  $\theta_1$ !



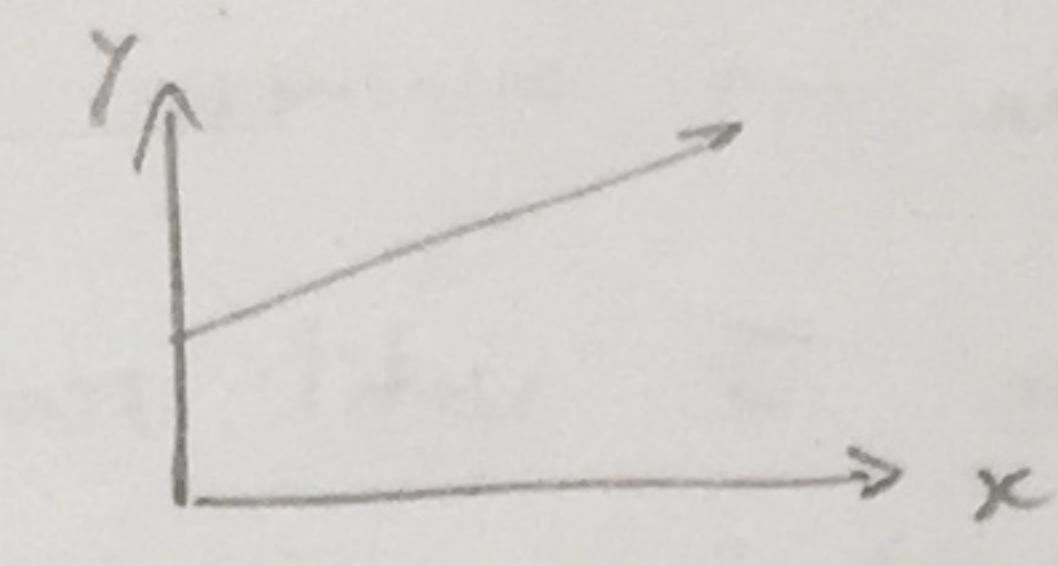
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

$$\theta_1 = 0.5$$



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

↳ Various  $\theta_0, \theta_1$ 's

⇒ want  $\theta_0$  &  $\theta_1$  that fits data! →  $h_{\theta}(x)$  close to  $y$  for  $(x, y)$  (training examples)

↳ minimize  $h_{\theta}(x) - y$  → square it → minimize  $(h_{\theta}(x) - y)^2$

over  $\theta_0, \theta_1$

↳ minimize  $\frac{1}{2}$  average sum-of-squares error

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

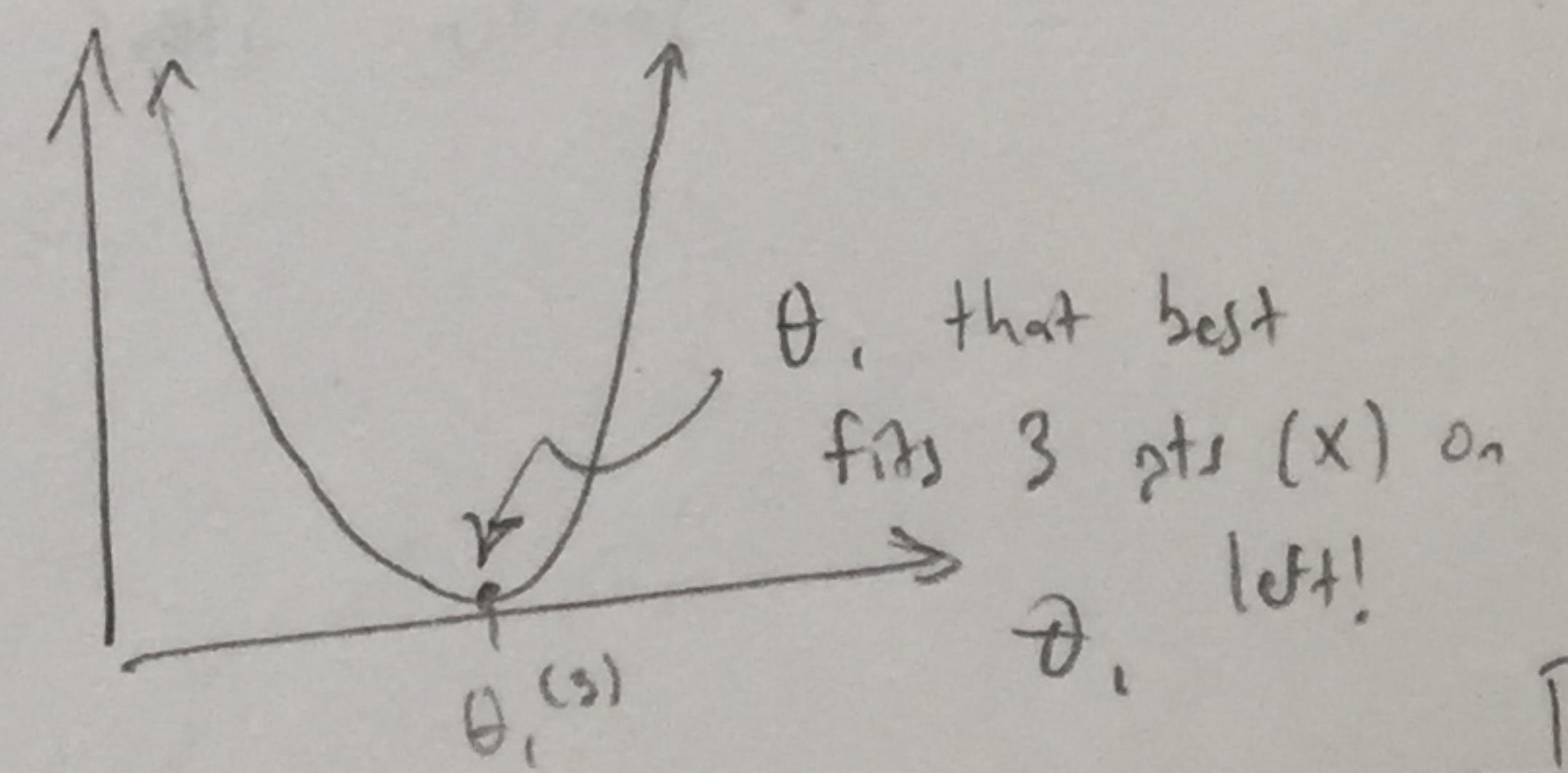
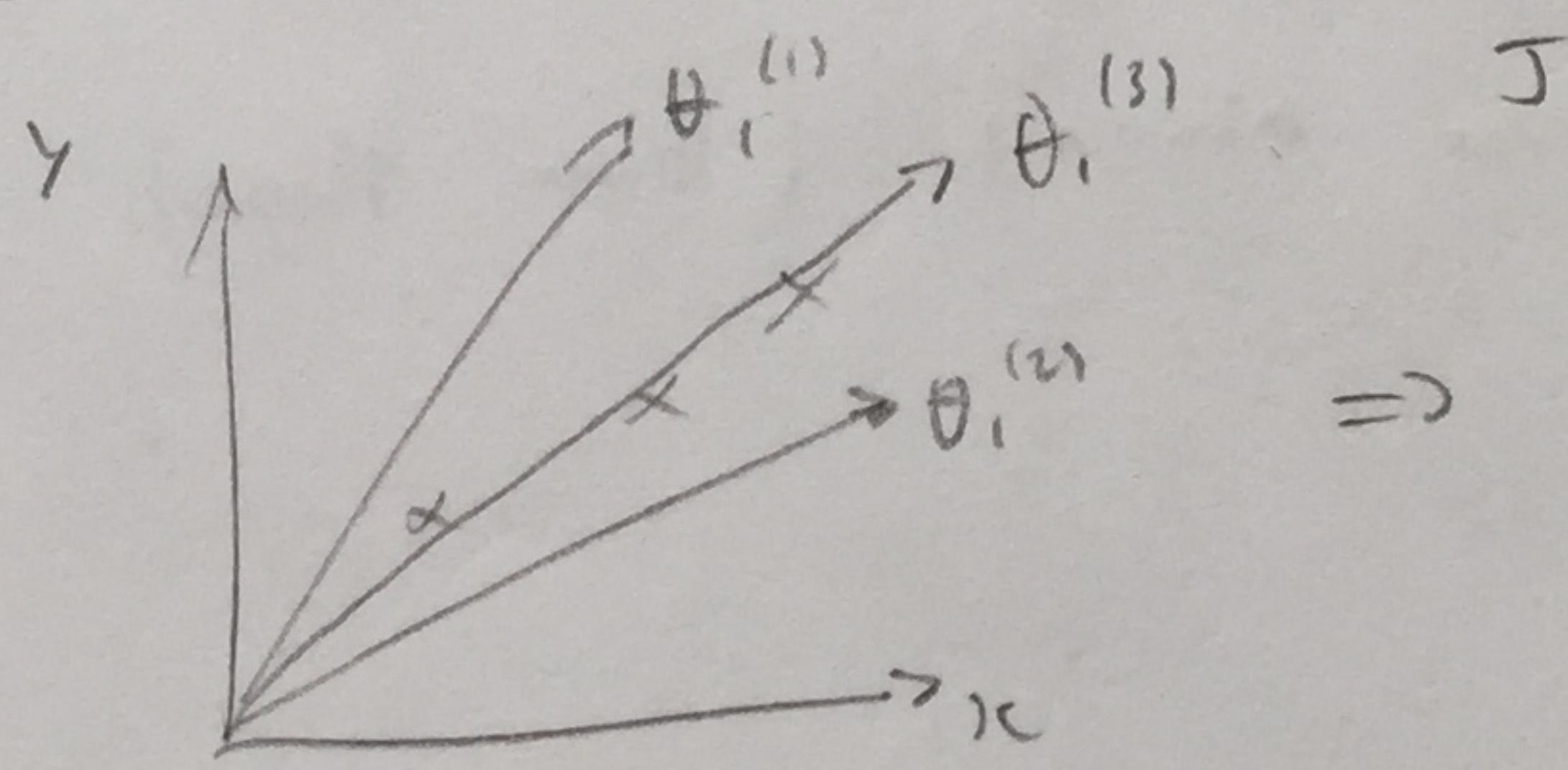
↳ minimize  $J(\theta_0, \theta_1)$  "squared error function"

⇒ minimize  $J(\theta_0, \theta_1) = \frac{1}{2} \bar{s}^2$ , where  $\bar{s}$  is mean of squares  $h_{\theta}(x^{(i)}) - y^{(i)}$

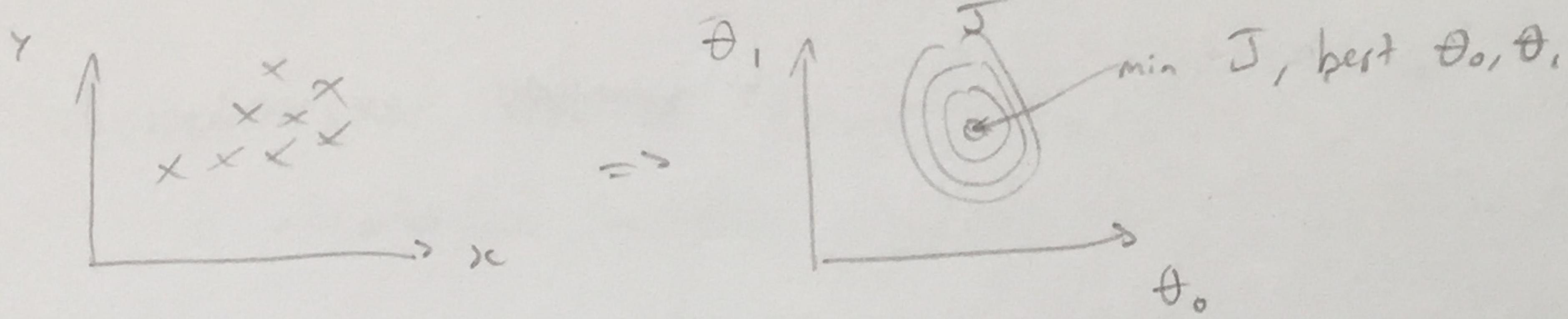
↳  $h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$

↳  $\frac{1}{2}$  makes derivative easier for gradient descent later

⇒ 2 plots:



- If  $\theta_0$  not fixed,  $J(\theta_0, \theta_1)$  plotted on contour / 3D plot!



- Gradient Descent:

↳ a minimization algorithm  $\rightarrow$  minimize  $J(\theta_0, \theta_1, \theta_2, \dots)$

↳ change  $\theta_i$  to reduce  $J$  until reach (hopefully) minimum

↳ initial choice matters!  $\rightarrow$  can reach different local optimums depending on initial choice

Algorithm:  $\text{temp } 0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$\alpha$  = step size

$$\text{temp } 1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

("learning rate")

$$\theta_0 := \text{temp } 0$$

$$\theta_1 := \text{temp } 1$$

$\Rightarrow$  Note order of operations!

"Simultaneous Update"

$\hookrightarrow :=$  means assignment ( $a := b$  means val.  $a = \text{val. } b$ )

ex)  $\theta_0 = 1$      $\theta_1 = 2$      $\theta_j := \theta_j + \sqrt{\theta_0 \theta_1}$   $\rightarrow$  find next  $\theta_0$  &  $\theta_1$

$$\Rightarrow \text{temp } 0 := 1 + \sqrt{1 \cdot 2} = 1 + \sqrt{2}$$

$$\text{temp } 1 := 2 + \sqrt{1 \cdot 2} = 2 + \sqrt{2}$$

$$\boxed{\theta_0 = 1 + \sqrt{2}}$$

$$\boxed{\theta_1 = 2 + \sqrt{2}}$$

$\Rightarrow \alpha$  too small  $\rightarrow$  SLOW, but  $\alpha$  too large  $\rightarrow$  possibly overshoot!

$\Rightarrow \alpha \cdot \frac{\partial}{\partial \theta_i} J(\theta_0, \dots)$  dictates  $\theta$  motion  $\rightarrow$  both step size & derivative matter

↳ smaller steps near minimum (lower slope)

## Regression w/ One Variable:

recall → linear regression model:  $h_{\theta}(x) = \theta_0 + \theta_1 x$  hypothesis (line here)

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

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

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

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

⇒ we need partial deriv term:  $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \left[ \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 \right]$

↳ for  $j=0$ : part. deriv =  $\frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) \cdot (1-0)$

for  $j=1$ : part. deriv =  $\frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) \cdot (x^{(i)} - 0)$

↳ gradient descent becomes:  $\text{temp } 0 := \theta_0 - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) \right]$   
 $\text{temp } 1 := \theta_1 - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x^{(i)} \right]$

$$\begin{aligned} \theta_0 &:= \text{temp } 0 \\ \theta_1 &:= \text{temp } 1 \end{aligned} \quad \text{simultaneous update}$$

⇒ for linear regression always bowl-shaped (convex function)

↳ only 1 global optimum → great for gradient descent!

⇒ "batch" gradient descent → each step uses all training examples  $\left( \sum_{i=1}^m \right)$