# Gems of TCS

## Linear Regression

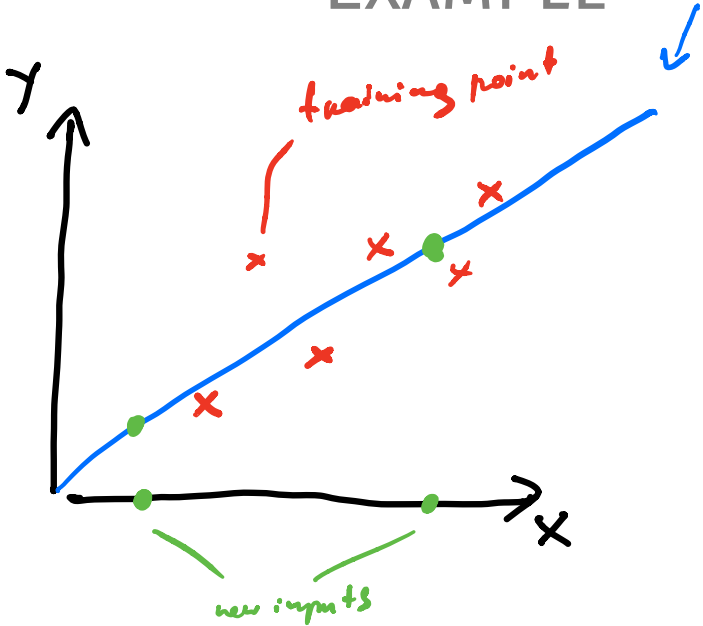Sasha Golovnev

May 6, 2021

# CLASSES OF LEARNING PROBLEMS
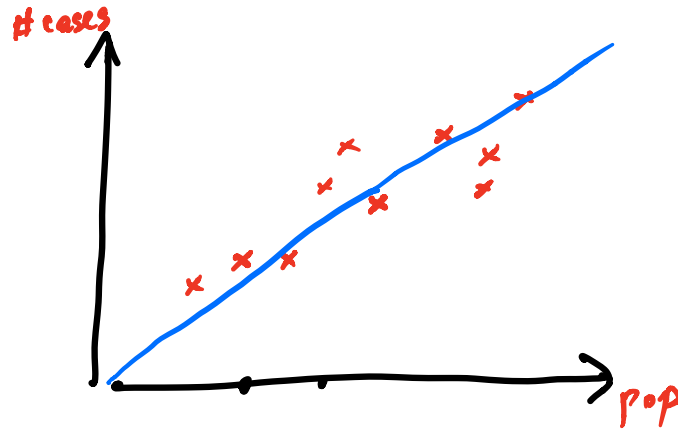
0   non-spam
1   spam
2   phishing

- Classification

- Ranking

- Regression      $\rightarrow$ real-valued output

- Clustering

- …

EXAMPLE
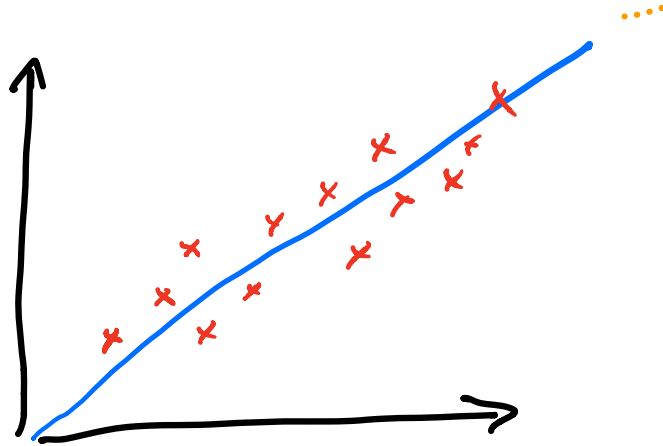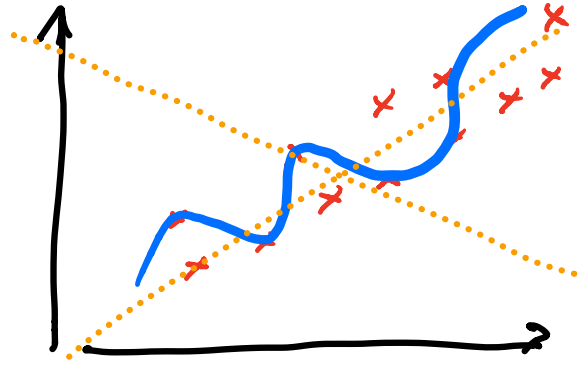
# Example

- Predict # of covid cases given city population

# LINEAR REGRESSION

predictor (hypothesis) is a linear function

$X_1 \qquad Y_1$

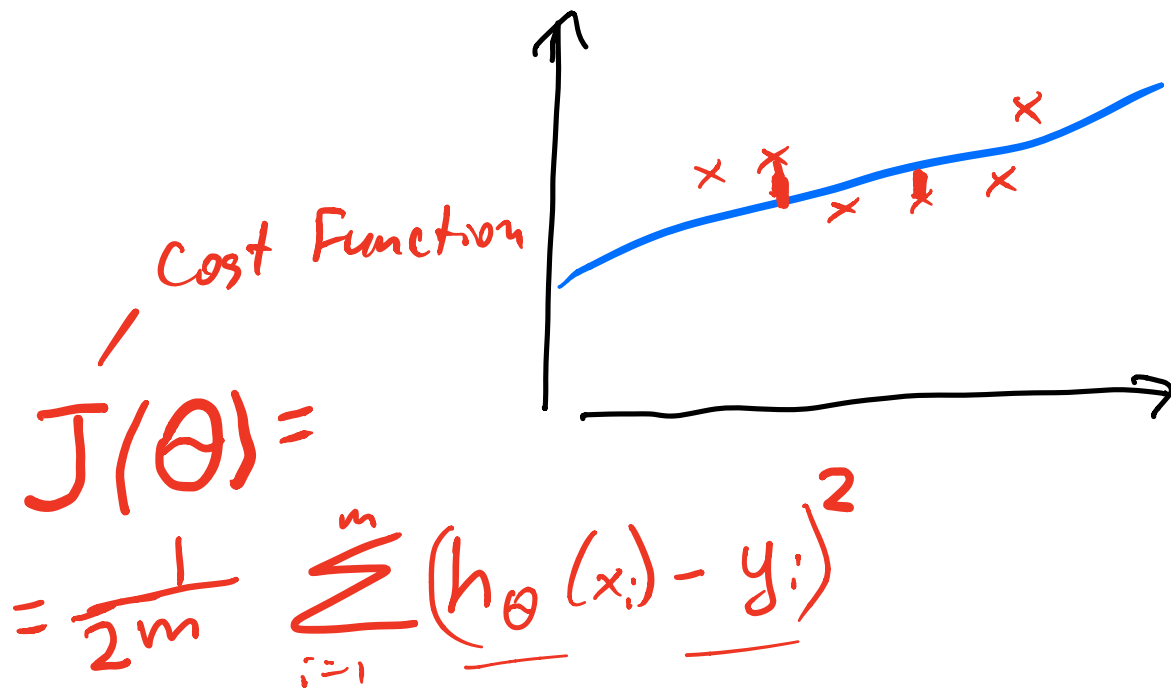$X_2 \qquad Y_2$

$\vdots$

$X_m \qquad Y_m$

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

$$\theta = (\theta_0, \theta_1)$$

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

Cost Function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x_i) - y_i \right)^2$$

# Linear Regression

Parameters of $(\theta_0, \theta_1)$

Hypothesis

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

Cost function:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x_i) - y_i)^2$$
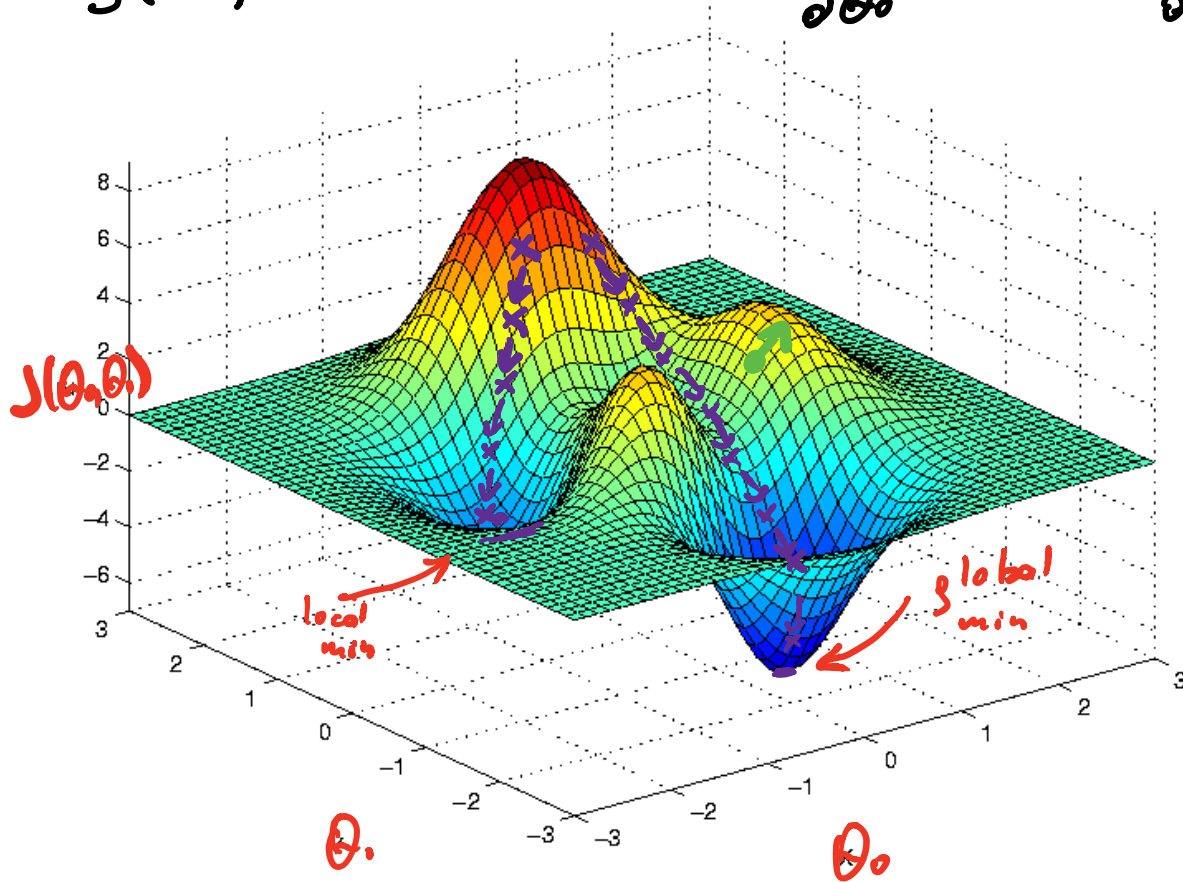
Goal:

minimize $J(\theta)$

Find $\theta = (\theta_0, \theta_1)$

# GRADIENT DESCENT

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

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} \qquad \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$



$J(\theta_0, \theta_1)$

local min
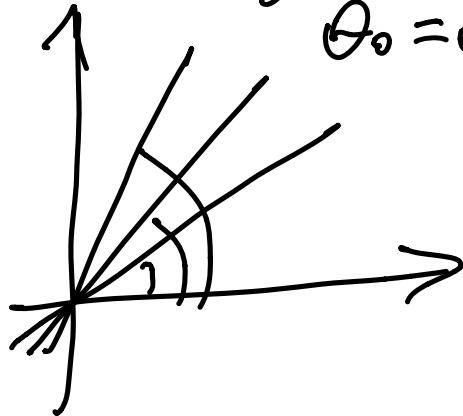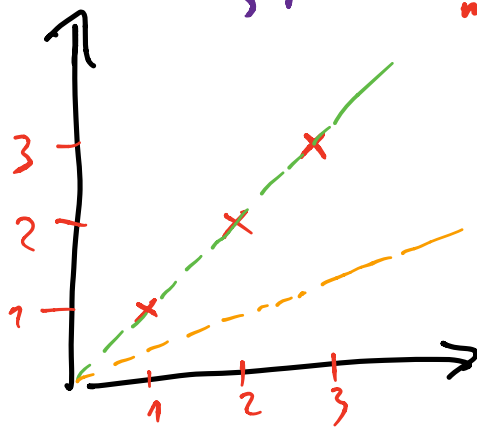
global min

$\theta_1$

$\theta_0$

# Toy Example

consider lines pass $(0,0)$

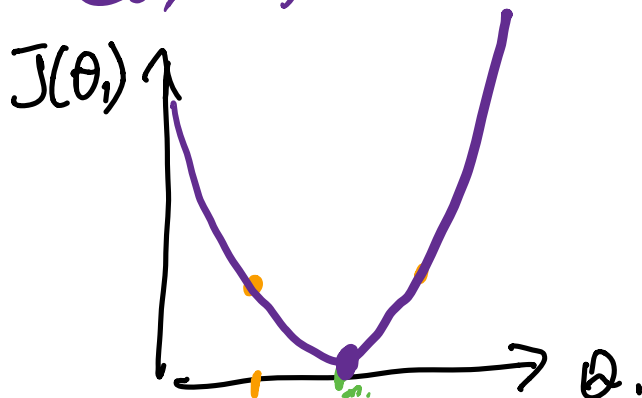$y = \theta_1 \cdot x$

$\theta_0 = 0$

## Training points

$m = 3$

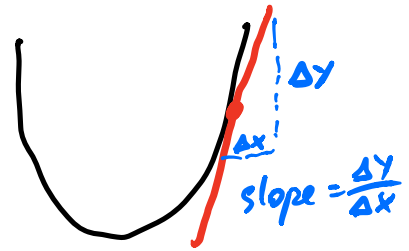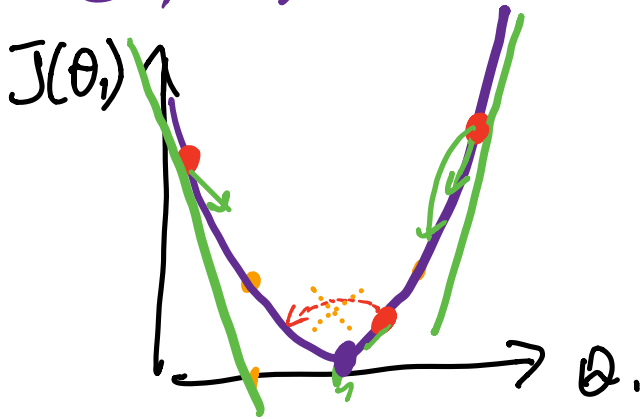$$J(\theta_1) = \frac{1}{6} \sum_{i=1}^{3} \left( \theta_1 x_i - y_i \right)^2$$

$\theta_1 = 1 \quad J(\theta_1) = 0$

$\theta_1 = 0.5 \quad J(\theta_1) =$

$= \frac{1}{6} \left( \left(\frac{1}{2}\right)^2 + 1^2 + \left(\frac{3}{2}\right)^2 \right) =$

$0.58$

## Cost function

$J(\theta_1)$

$\theta_1$

# Cost function



- function increases (decreases), we want to decrease (increase) $\theta_1$

- function increases rapidly, we want to decrease $\theta_1$ a lot

  slowly, we want to decrease $\theta$ a bit

$J'(\theta_1)$ is positive $\Leftrightarrow$ $J$ increases at $\theta_1$

negative $\Leftrightarrow$ $J$ decreases at $\theta_1$

increases slowly $\Leftrightarrow$ derivative is small

fast $\Leftrightarrow$ large

$$\theta_1 = \theta_1 - \boxed{\alpha} \cdot J'(\theta_1)$$

Learning rate

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

increases fast $\left( \dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}, \right.$

$$\left. \dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \right)$$

decreases fast

$$\left( -\dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}, -\dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \right)$$

# Gradient Descent

Pick $\theta_0, \theta_1$ (Say, $\theta_0 = \theta_1 = 0$)

repeat until converge:

==simultaneously== update

$$\begin{cases} \theta_0 = \theta_0 - \boxed{\alpha \cdot \dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}} \\ \theta_1 = \theta_1 - \alpha \cdot \dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} \end{cases}$$
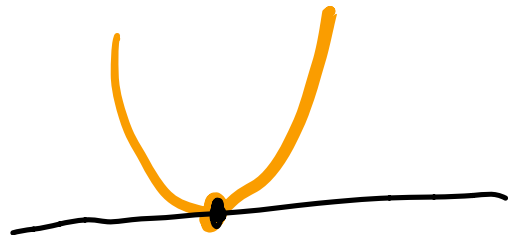
repeat until converge:

$$update_0 = \boxed{\alpha \cdot} \dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_0}$$

$$update_1 = \alpha \cdot \dfrac{\partial J(\theta_0, \theta_1)}{\partial \theta_1}$$

$$\theta_0 = \theta_0 - update_0$$

$$\theta_1 = \theta_1 - update_1$$

learning rate

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (\theta_0 + \theta_1 x_i - y_i)^2$$

$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^{m} (\theta_0 + \theta_1 x_i - y_i)$$
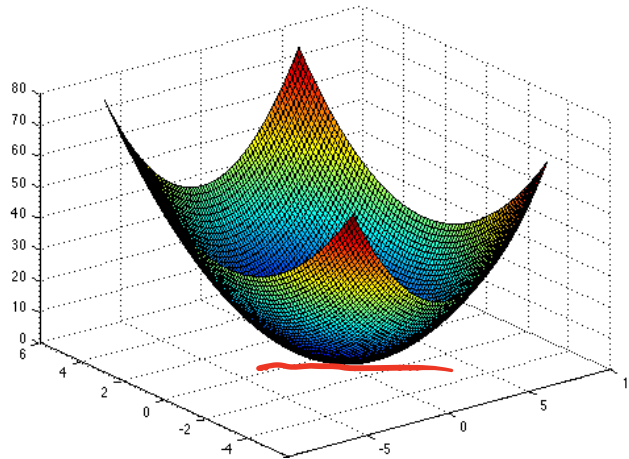
$$\frac{\partial J(\theta_0, \theta_1)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^{m} (\theta_0 + \theta_1 x_i - y_i) \cdot x_i$$

$$\varepsilon = \frac{1}{10^3}$$

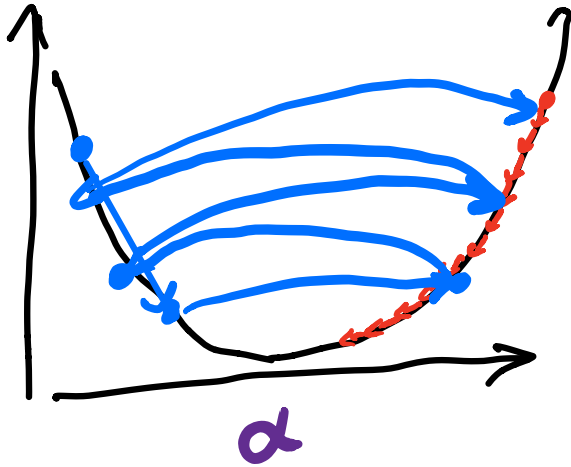in practice, converge $\equiv$ update $< \varepsilon$

Cost function for linear regression has one local minimum

# LEARNING RATE

Learning rate $\alpha$

Thm For small enough, gradient descent will converge.

too small $\alpha$ — slow alg.
too big $\alpha$ — doesn't converge



$\alpha$
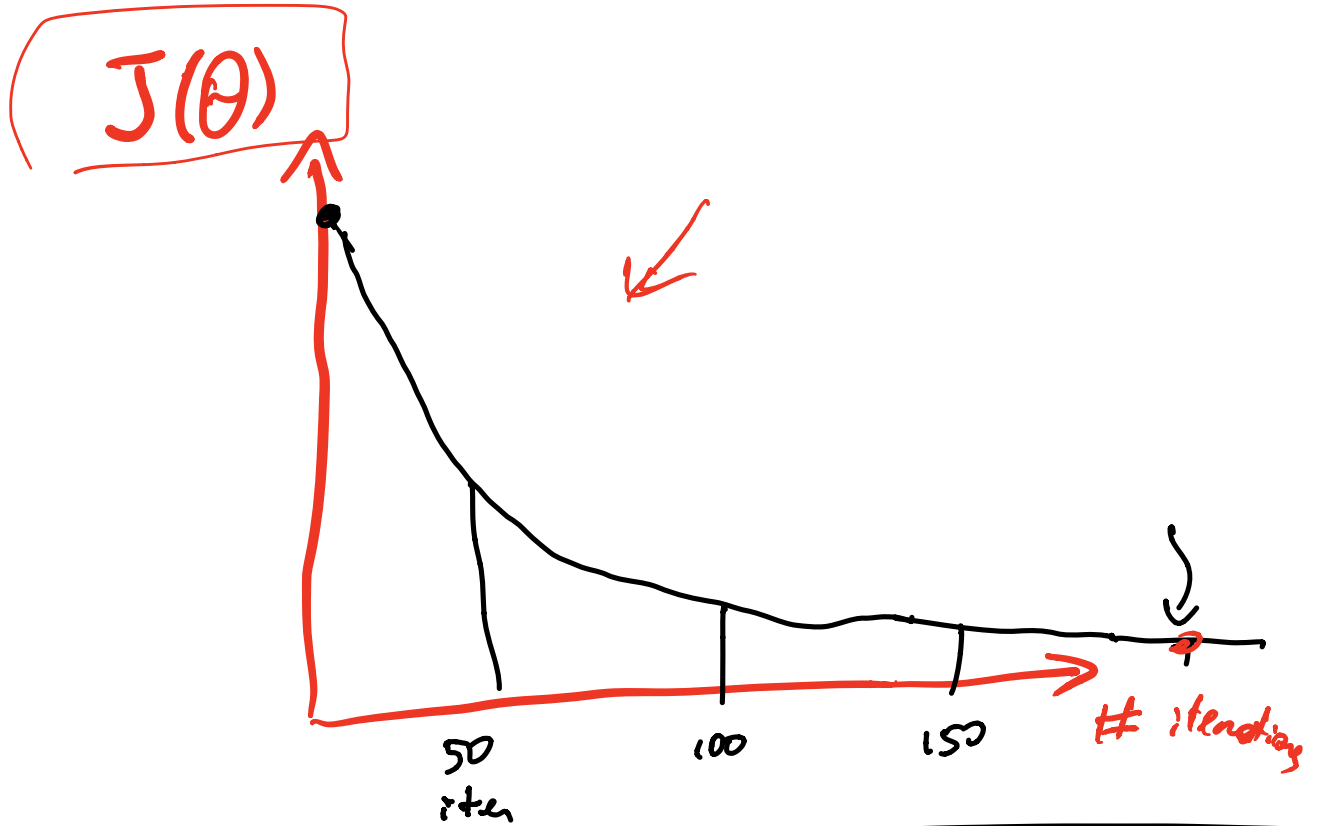
--- $\frac{1}{1000}$ $\frac{1}{100}$ $\left(\frac{1}{10} \ \frac{1}{5} \ \frac{1}{3} \ \frac{1}{2} \ 1\right)$ 10 100 1000 ---

slow slow

doesn't
conv

# "Debug" Gradient Descent

$J(\theta)$

# iterations

50 iter

100

150

$\alpha$ is too big