

# ARTIFICIAL INTELLIGENCE

LEI/3, LMA/3, MBE/1

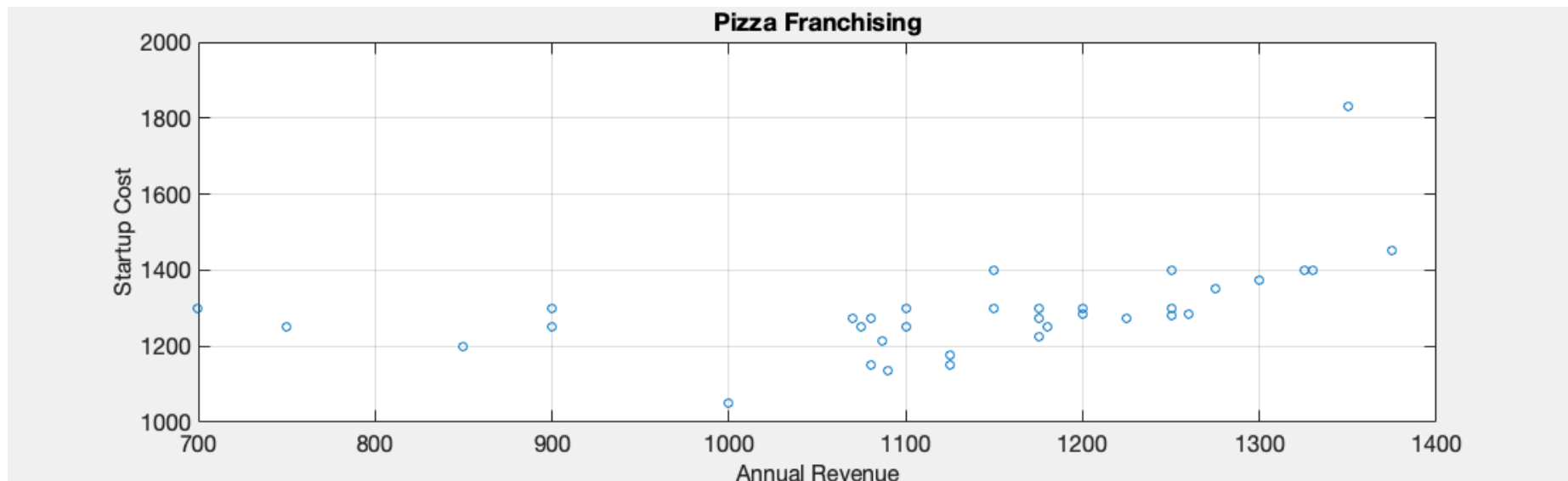
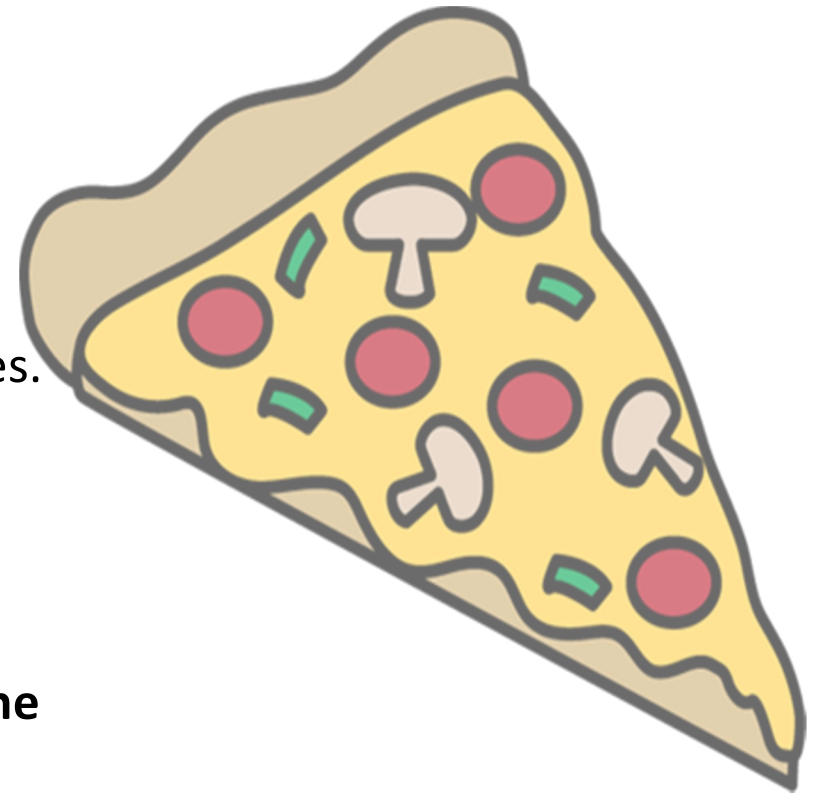
University of Beira Interior, Department of Informatics

Hugo Pedro Proença

[hugomcp@di.ubi.pt](mailto:hugomcp@di.ubi.pt), 2022/23

# Pizza Franchising

- Pizza is a \$45.1 billion industry in the United States.
- Suppose that one of the most well-known Pizza chain is interested in perceiving the relationship between the **average annual revenue** of its local stores and the corresponding **startup cost**.
- This data will be of maximum interest to **define the franchise fee for future openings**



# Pizza Franchising

- It appears that there is a **direct relation** between the annual income of one store, and the cost to start the store.
  - On average, larger stores sell more Pizza, but also they are more costly to set up:
  - Furniture, taxes, employees...
- In this problem we have 36 examples, typically designated as “**instances**”
  - $N=36$
- The independent variables are typically referred to as “**features**”
  - Are the input variables ( $x$ )
- The number of features determines the **dimensionality of the problem**
  - $d=1$
- The dependent variable is typically designated as the output, or “**target**”
  - The target distribution determines the type of supervised machine learning problem: classification or **regression** (in this case)

Annual Fee	Startup Cost
1000	1050
1125	1150
1087	1213
1070	1275
1100	1300
1150	1300
1250	1400
1150	1400
1100	1250
1350	1830
1275	1350
1375	1450
1175	1300
1200	1300
1175	1275
1300	1375
1260	1285
1330	1400
1325	1400
1200	1285
1225	1275
1090	1135
1075	1250
1080	1275
1080	1150
1180	1250
1225	1275
1175	1225
1250	1280
1250	1300
750	1250
1125	1175
700	1300
900	1250
900	1300
850	1200

Independent Variable

Dependent Variable

# Machine Learning I: Model Representation

- Suppose that the experts/administration/managers of the Pizza chain think that it might exist a roughly linear relationship between the annual revenue of one store and its startup cost:
  - This kind of “expertise” is always valuable to machine learning, as it simplifies the range of models that we can attempt to create
- Also, one of the Machine Learning’s foundation is the **Occam’s razor**:
- Known as the law of parsimony
  - Is a problem-solving principle that essentially states that "**simpler solutions are more likely to be correct than complex ones**".
  - When comparing competing hypotheses to solve a problem, one should select the solution with the fewest assumptions, i.e., **the simplest**
- The idea is attributed to English Franciscan friar William of Ockham (1287–1347), a scholastic philosopher and theologian.

# Machine Learning I: Model Representation

- Linear Model
  - According to Occam's razor (and the administration also!), in the Pizza Franchising, we should start by consider a purely linear model to “**describe the pattern**” (i.e., describe the relationship) between the independent(s) and the dependent variables
- Formally, our model (**hypothesis**) is that:

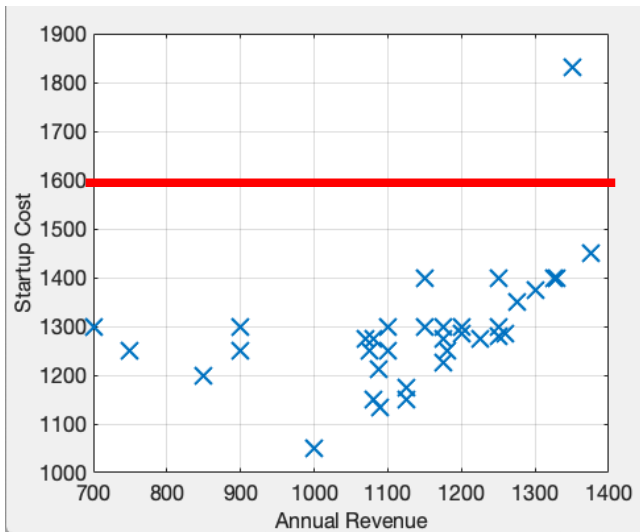
$$h_{\theta}(x) = \theta_1 \cdot x + \theta_2$$

- The task of Machine Learning is to **find us the best possible model**, i.e., the one **that optimally expresses the relationship** between the independents and dependent variables
- This essentially involves to find the optimal  $(\theta_1, \theta_2)$  values
- After all, we end up with an **optimization problem in the  $R^2$  space**

# Machine Learning II: Cost Function

- Clearly, there will be models that are better than others:

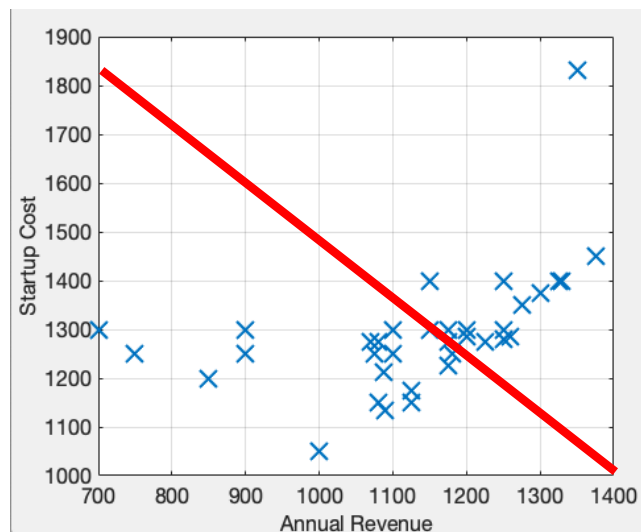
$$(\theta_1 = 0, \theta_2 = 1600)$$



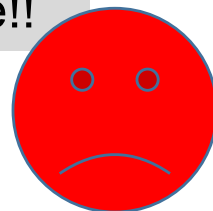
Bad.



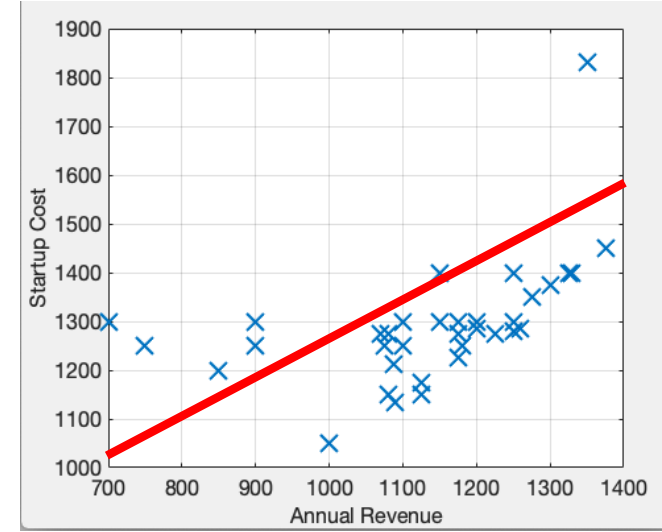
$$(\theta_1 = -1.15, \theta_2 = 1005)$$



Terrible!!



$$(\theta_1 = 0.82, \theta_2 = 446)$$



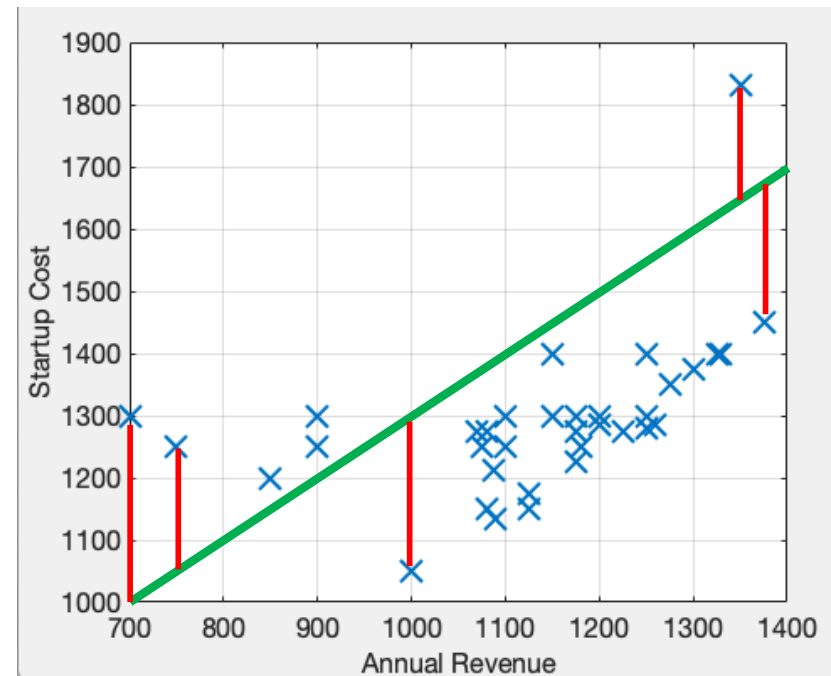
Good...



...but ***“the best”***?

# Machine Learning II: Cost Function

- The **Cost Function** should distinguish between two alternate hypotheses, i.e., it should be used to favor one hypothesis instead of other
- In practice, the cost function receives the **parameters of one model** and returns “*how good/bad the model is*”
- In this problem, we are interested in models that are as close as possible to the data points
- I.e., the “optimal model” will overlap exactly all the points we have in the dataset
  - Impossible, for the type of model chosen



# Machine Learning II: Cost Function

- The **Cost Function** is typically expressed as **J()**
- **The cost function receives as input, the parameters of the model**
  - In this case, it receives two parameters:  $(\theta_1, \theta_2)$
- Hence, the cost function is formally  $J: \mathbb{R}^2 \rightarrow \mathbb{R}$

$$J(\theta_1, \theta_2) = \frac{1}{2N} \sum_{i=1}^N (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

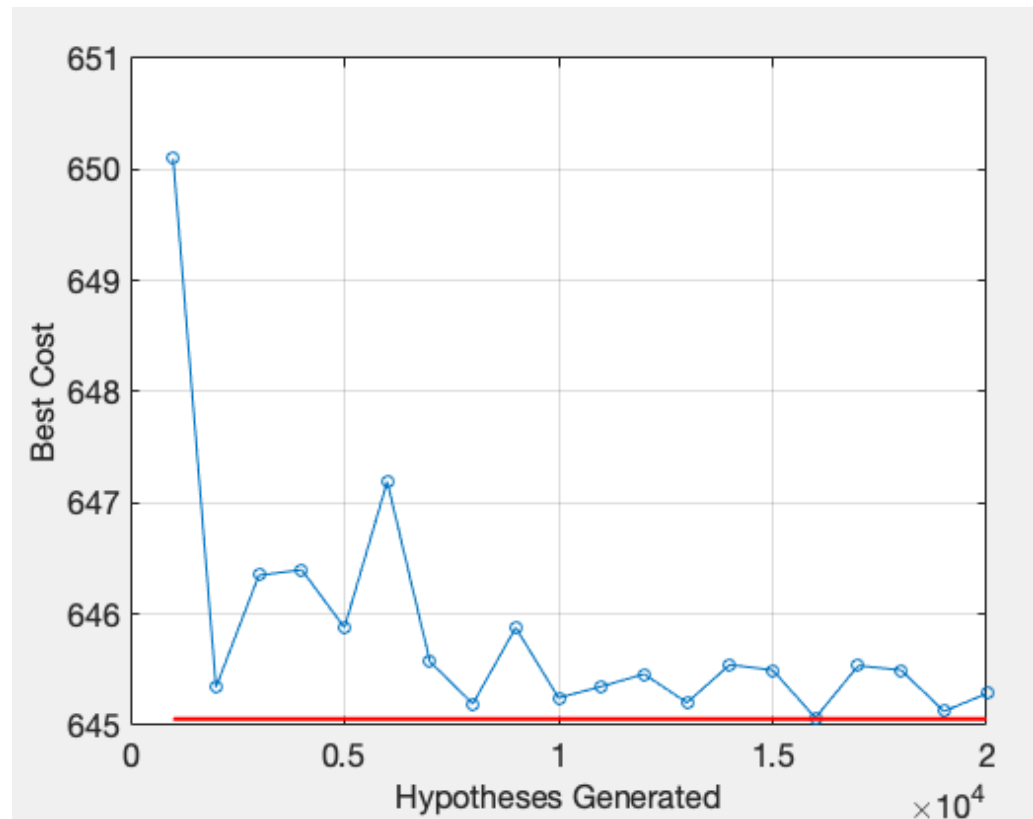
Why ??

- In practice, this function sums up all the Euclidean distances between the targets (ground truth) in our dataset and the values given by the model at each point
  - Clearly, if one model is optimal  $h_{\theta}(x^{(i)}) == y^{(i)}$  and **J=0**
- At the (almost) end of this story, Machine Learning is about minimizing **J()**



# Machine Learning III: Optimization

- “Computers are so fast these days, what if we simply generate millions of different hypotheses and pick the **best one?**”
  - This is the “brute-force” approach, that (only) in problems of reduced dimensionality might lead to reasonable results.
- The plot given at right compares the best model obtained “*by chance*” (dependent variable), with respect to the numbers of models randomly created (independent variable).
- In some cases, the best random model was “close” to the optimal model:
  - Cost 645.05
  - $(\theta_1, \theta_2) = (0.376, 867.6)$



# Machine Learning III: Optimization

- How to obtain the best possible model?
- Find the  $(\theta_1, \theta_2)$  parameters that minimize  $J()$
- Formally:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} J(\theta_1, \theta_2)$$

- In practice, this is an optimization problem in 2D space, that requires to find the derivative of  $J()$  with respect to  $\boldsymbol{\theta}$ .
- Recall from single variable calculus that (assuming a function  $f$  is differentiable) the minimum  $x^*$  of  $f$  has the property that the derivative  $df/dx$  is zero at  $x=x^*$ 
  - An analogous result holds in the multivariate case:

The diagram illustrates the decomposition of the gradient of the cost function  $J(\boldsymbol{\theta})$  into its partial derivatives. On the left, the expression  $\frac{\partial}{\partial \boldsymbol{\theta}} J(\boldsymbol{\theta})$  is shown. Two arrows originate from this expression and point to the right, where the two partial derivatives are listed:  $\frac{\partial}{\partial \theta_1} J(\boldsymbol{\theta})$  (top) and  $\frac{\partial}{\partial \theta_2} J(\boldsymbol{\theta})$  (bottom).

Partial  
Derivatives

# Machine Learning Optimization: Closed-Form

- Minimizing  $J()$  is equivalent to minimize:

$$\sum_{i=1}^N (\theta_1 x^{(i)} + \theta_2 - y^{(i)})^2$$

$$\mathbf{X} = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \end{bmatrix}$$

**Bias!!**

- Using matrix algebra, we know that

$$\sum_{i=1}^N (\theta_1 x^{(i)} + \theta_2 - y^{(i)})^2 = (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$$

- So, we are interested in minimizing the above expression, i.e.,

$$\frac{\partial}{\partial \boldsymbol{\theta}} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y}) = 0$$

- Applying the distributive property. Also:

$$(\mathbf{AB})^\top = \mathbf{A}^\top \mathbf{B}^\top$$

$r$  scalar,  $r^\top = r$   
 $\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta}$  is scalar.

$$\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{X}^\top \boldsymbol{\theta}^\top \mathbf{X}\boldsymbol{\theta} - \mathbf{X}^\top \boldsymbol{\theta}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} + \mathbf{y}^\top \mathbf{y} = 0$$

$$\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta} = (\mathbf{y}^\top \mathbf{X}\boldsymbol{\theta})^\top = \mathbf{y} \mathbf{X}^\top \boldsymbol{\theta}^\top$$

# Machine Learning Optimization: Closed-Form

- Simplifying:

$$\frac{\partial}{\partial \theta} \mathbf{X}^T \theta^T \mathbf{X} \theta - 2 * \mathbf{X}^T \theta^T \mathbf{y} + \mathbf{y}^T \mathbf{y} = 0$$

**Matrix Derivatives:**

$$\frac{\partial (\mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \mathbf{A}^T$$

$$\frac{\partial (\mathbf{X}^T \mathbf{X})}{\partial \mathbf{X}} = 2\mathbf{X}$$

- Applying the derivatives rules:

$$\frac{\partial (\mathbf{X}^T \mathbf{A}\mathbf{X})}{\partial \mathbf{X}} = \mathbf{A}\mathbf{X} + \mathbf{A}^T \mathbf{X}$$

$$2 \mathbf{X}^T \mathbf{X} \theta - 2 * \mathbf{X}^T \mathbf{y} = 0$$

$$\begin{aligned} \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{y} &= 0 \\ \mathbf{X}^T \mathbf{X} \theta &= \mathbf{X}^T \mathbf{y} \\ (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) \theta &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

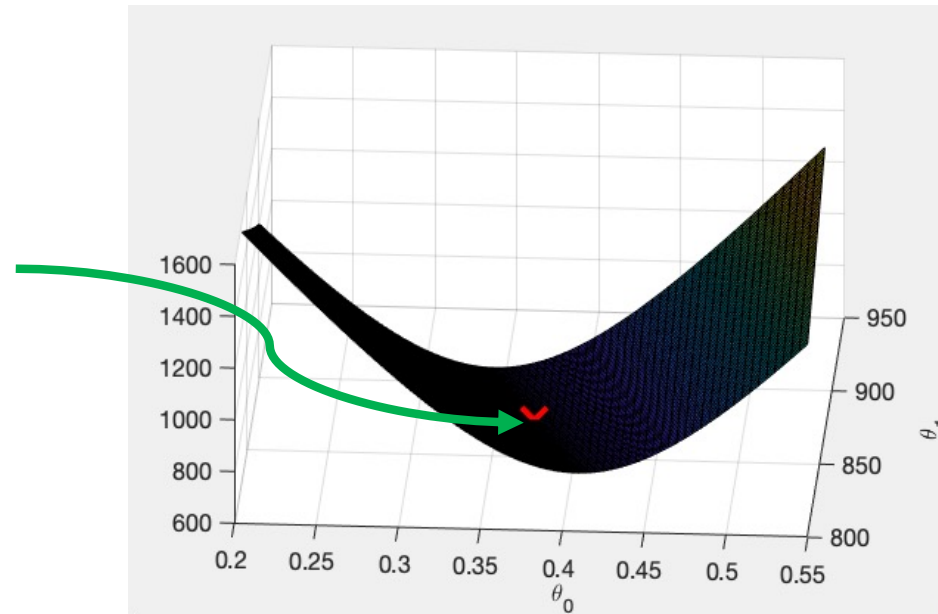
- Solving with respect to  $\theta$  :

$$\theta^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

# Machine Learning Optimization: Closed-Form

- The closed-form solution should be preferred for “smaller” datasets
  - When computing the matrix inverse is not a concern.
- For very large datasets, obtaining  $(\mathbf{X}^T\mathbf{X})^{-1}$  can be extremely costly
  - $\mathbf{X}$  has  $N \times (d+1)$  dimensions
- Also, there are cases where the  $(\mathbf{X}^T\mathbf{X})^{-1}$  not exists
  - e.g., the matrix is non-invertible (singular) in case of perfect multicollinearity

If succeeded, the Closed-Form enables us to obtain the optimal configuration of the hypothesis  $\theta^*$  in a single step



# Machine Learning Optimization: Partial Derivatives

- As we have seen, the goal is to obtain the  $\theta$  parameterization that minimizes  $J()$ :

$$J(\theta_1, \theta_2) = \frac{1}{2N} \sum_{i=1}^N (\theta_1 x^{(i)} + \theta_2 - y^{(i)})^2$$

- $(a+b)' = a' + b'$

$$\frac{\partial}{\partial \theta_1} J(\theta) = \frac{1}{2N} \sum_{i=1}^N 2 (\theta_1 x^{(i)} + \theta_2 - y^{(i)}) x^{(i)}$$

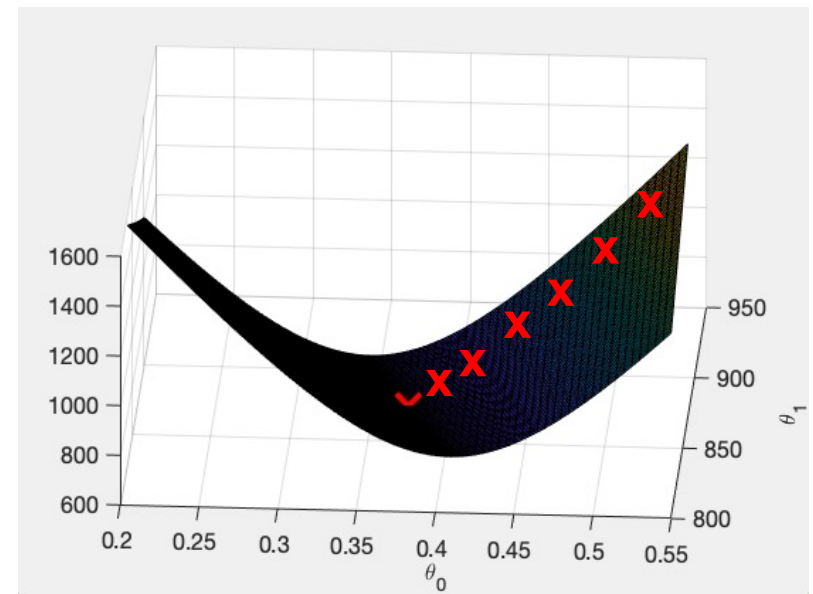
That's why!

$$\frac{\partial}{\partial \theta_2} J(\theta) = \frac{1}{2N} \sum_{i=1}^N 2 (\theta_1 x^{(i)} + \theta_2 - y^{(i)})$$

# Machine Learning Optimization: Gradient Descent

- In most practical cases, the Closed-Form is hard to obtain, and the solution is to use the “**Gradient Descent**” optimization version:
- Algorithm:
  1. Start with some random  $\theta$  configuration.  $\theta^{(0)}$
  2. Change iteratively (and slightly)  $\theta$  , to reduce  $J(\theta)$ 
    1.  $\theta^{(t+1)} = \theta^{(t)} - \Delta \frac{\partial}{\partial \theta} J(\theta)$
  3. (Hopefully) end up in a minimum

The rationale is to iteratively move in the steepest descend direction, in order to reach the (eventually local) minimum

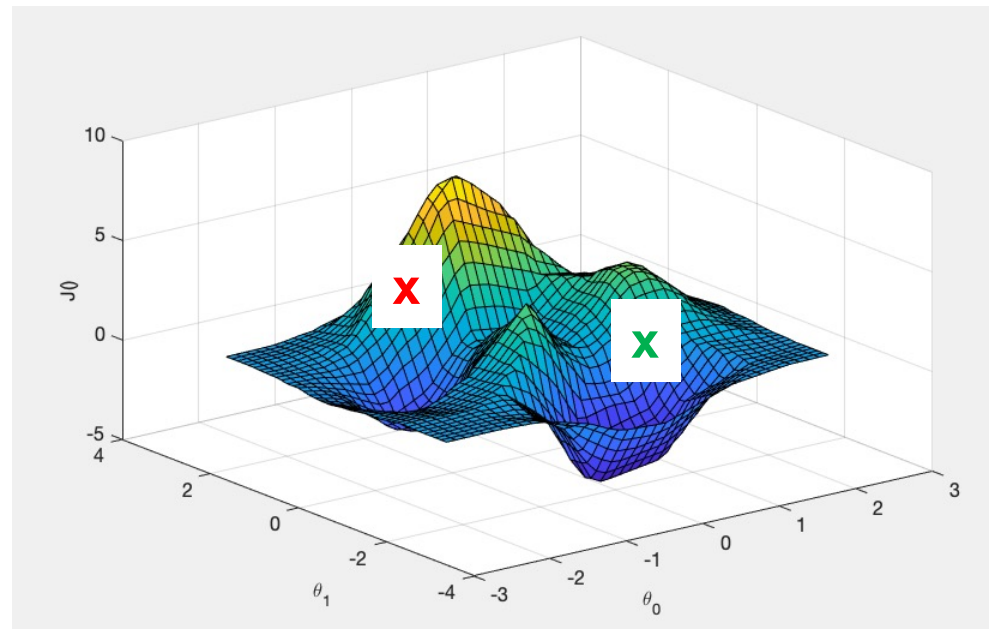


# Machine Learning Optimization: Gradient Descent

$$\theta_0 = \theta_0 - \Delta \frac{1}{N} \sum_{i=1}^N (\theta_1 x^{(i)} + \theta_2 - y^{(i)}) x^{(i)}$$

$$\theta_1 = \theta_1 - \Delta \frac{1}{N} \sum_{i=1}^N (\theta_1 x^{(i)} + \theta_2 - y^{(i)})$$

Main  
assumption in  
Gradient  
Descent:  
**Convexity!**





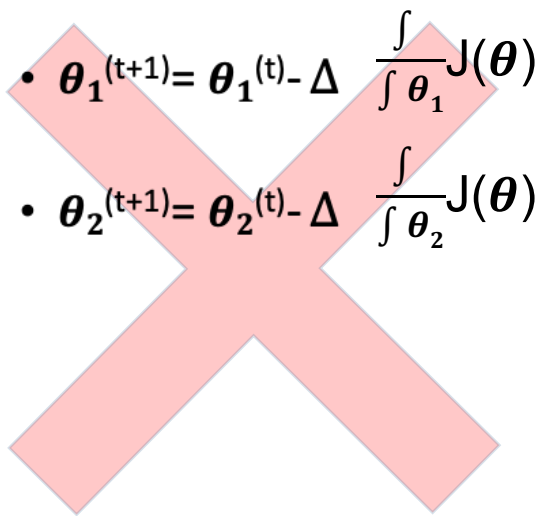
# Machine Learning Optimization: Gradient Descent

- Learning Rate

- Too large values lead to divergence
  - The optimal value of  $J()$  is not achieved, i.e., the best  $\theta$  configuration is not found
- Too small values slow down the learning process.

- Remark

- The update of parameters should be done simultaneously:



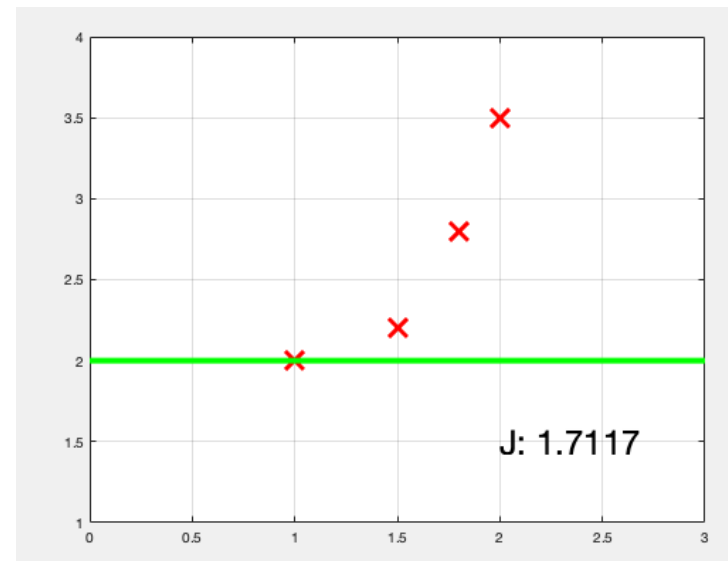
- $\theta_1^{(t+1)} = \theta_1^{(t)} - \Delta \frac{\partial}{\partial \theta_1} J(\theta)$
- $\theta_2^{(t+1)} = \theta_2^{(t)} - \Delta \frac{\partial}{\partial \theta_2} J(\theta)$

- $\text{aux}_1 = \theta_1^{(t)} - \Delta \frac{\partial}{\partial \theta_1} J(\theta)$
- $\text{aux}_2 = \theta_2^{(t)} - \Delta \frac{\partial}{\partial \theta_2} J(\theta)$
- $\theta_1^{(t+1)} = \text{aux}_1$
- $\theta_2^{(t+1)} = \text{aux}_2$

# Gradient Descent Exercise

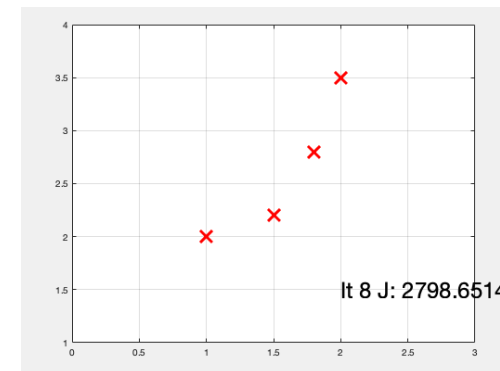
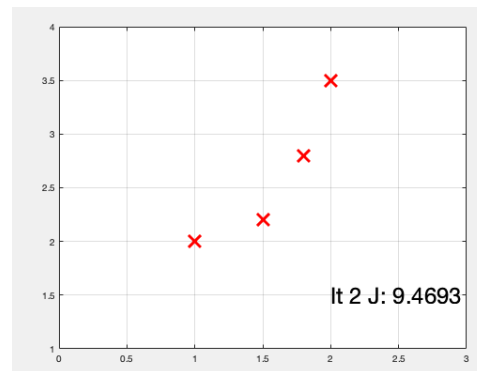
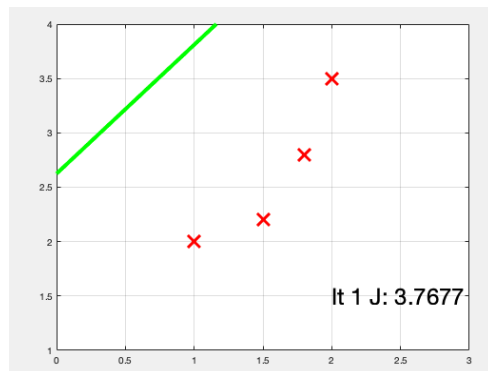
- Consider the following tiny dataset. Use the gradient descent algorithm to obtain the optimal linear regression hypothesis:
  - Start with  $\theta_1, \theta_2 = (0, 2)$
  - Use  $\Delta = 0.1$

X	Y
1	2
1.5	2.2
1.8	2.8
2	3.5



# Gradient Descent Exercise

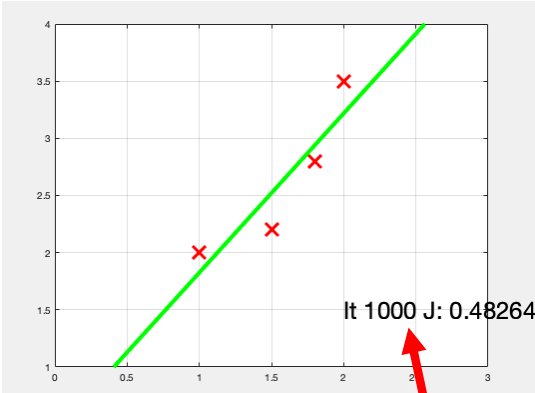
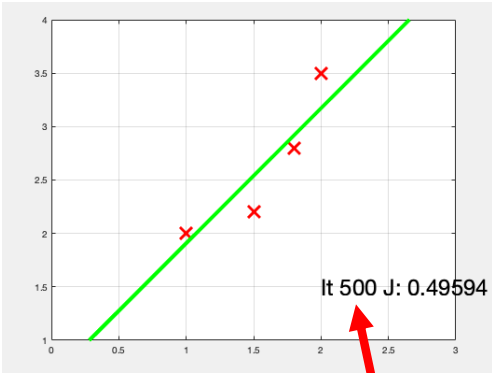
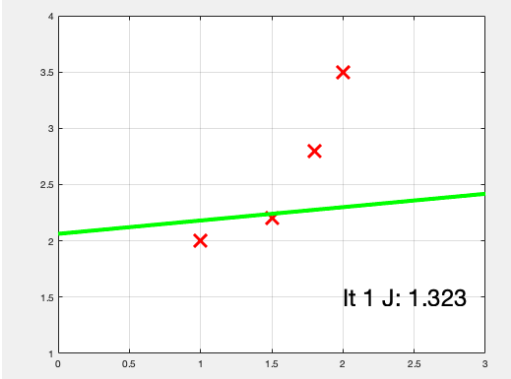
- Consider the following tiny dataset. Use the gradient descent algorithm to obtain the optimal linear regression hypothesis:
  - Start with  $\theta_1, \theta_2 = (0,2)$ 
    - Use  $\Delta=1$
    - Use  $\Delta=0.1$
    - Use  $\Delta=0.5$
- $\Delta=1$



**Diverged!!**

# Gradient Descent Exercise

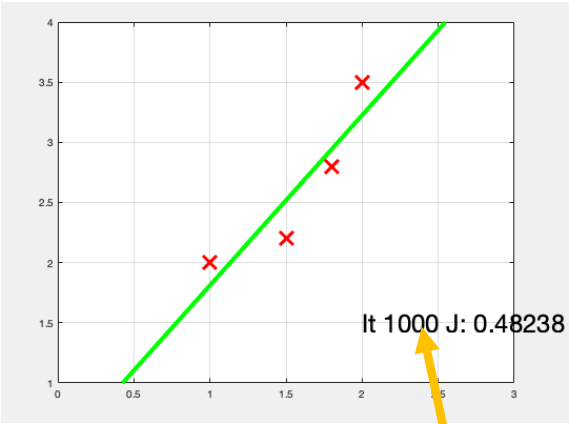
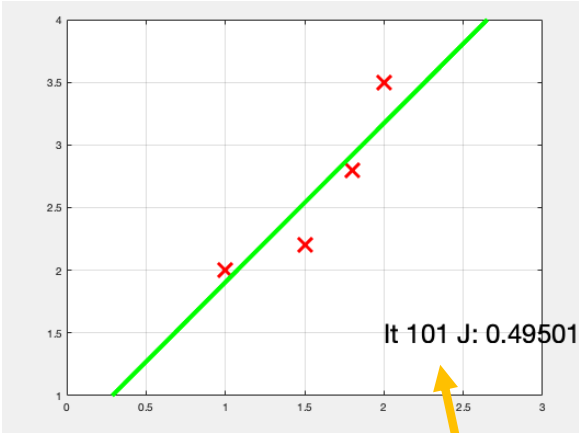
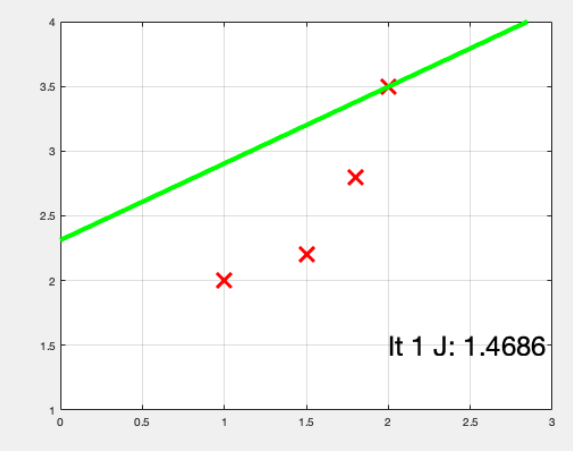
- $\Delta=0.1$



Too slow...

# Gradient Descent Exercise

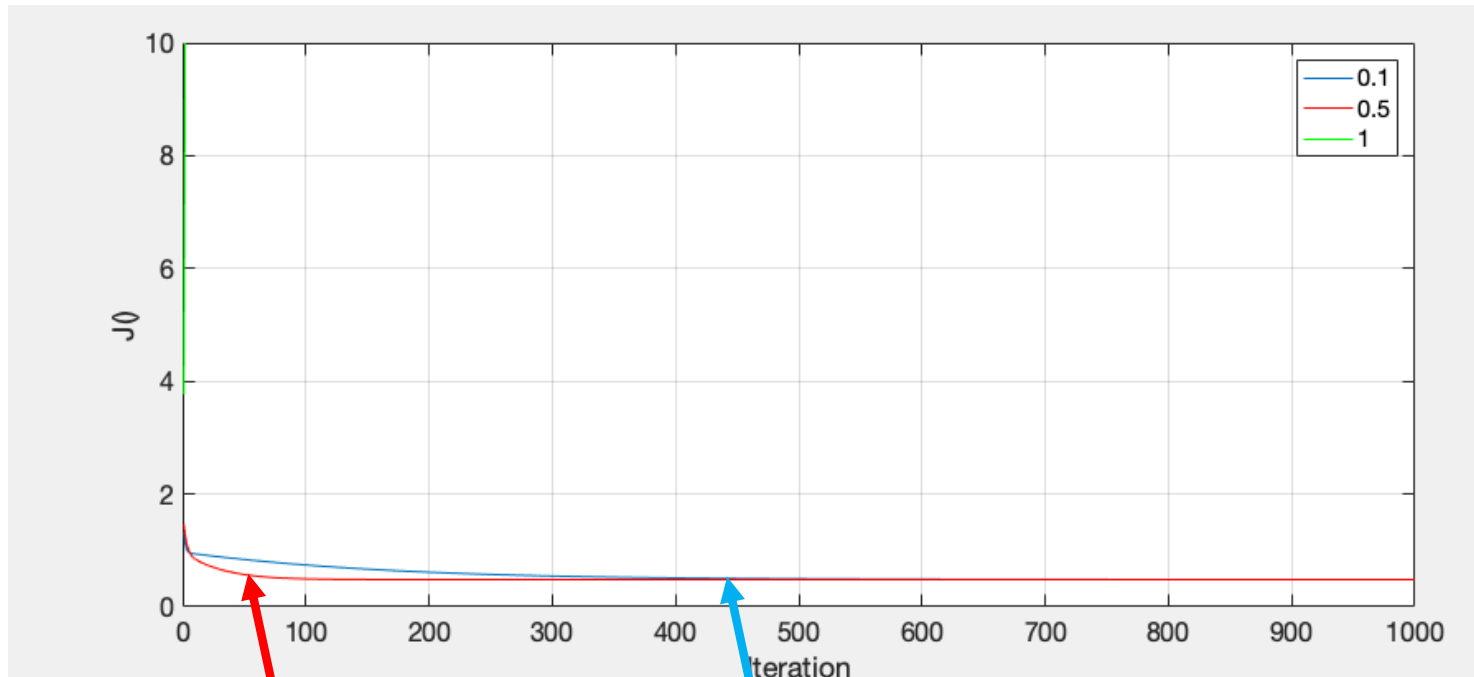
- $\Delta=0.5$



Better...

# Gradient Descent Exercise

- $\Delta=1$  vs.  $\Delta=0.1$  vs.  $\Delta=0.5$



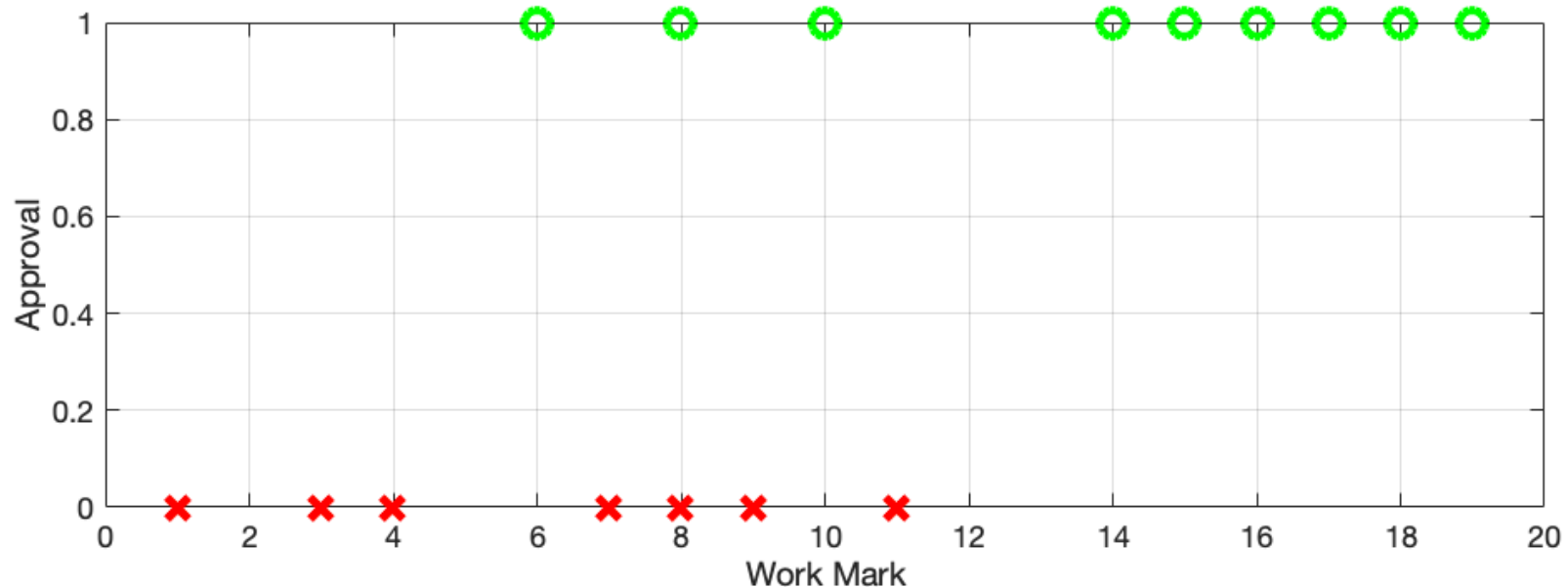
Acceptable solution ( $\Delta= 0.5$ )

Acceptable solution ( $\Delta= 0.1$ )

- Stop Criteria:
  - “T” iterations
  - While it stops to improv (i.e.,  $J^{(t+1)} - J^{(t)} < \varepsilon$  )

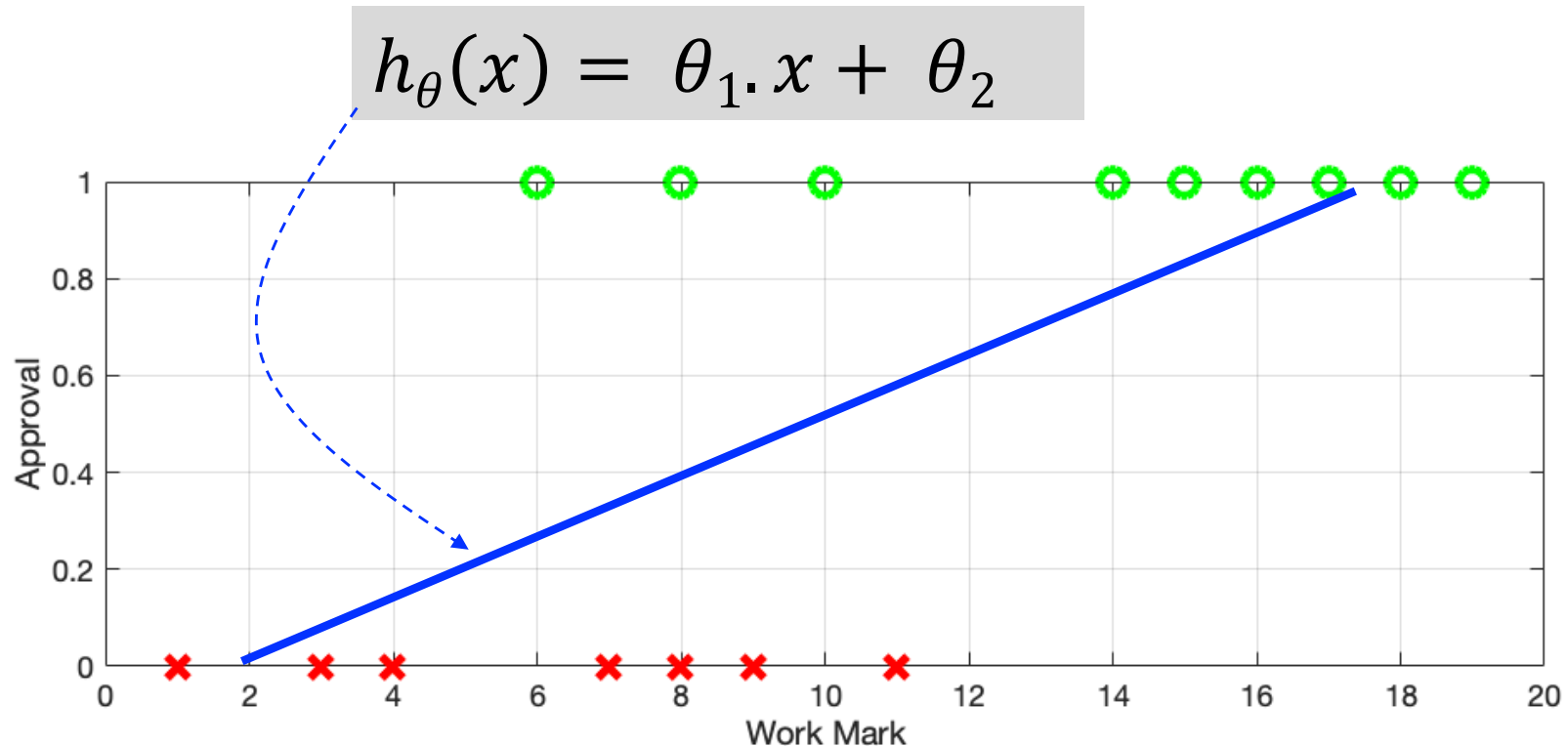
# Students Performance

- Suppose that we are interested in predicting the **approval rate** of a class, based on the **students marks in the first practical work**.
  - Typically, students that get good marks in the first work, got approved at the course.
  - Students with very low marks at the first work tend to fail in the final examination.
- Hence, our machine learning model is expected to predict a **binary outcome** (1: pass vs. 0: fail)



# Students Performance

- In this kind of problems, the dependent variable assumes a reduced set of labels:
  - Emails: “is this a **spam** or **no spam** email”?  $y \in \{0, 1\}$
  - Medical diagnosis: “is the patient **ill** or **healthy**”  $y \in \{0, 1\}$
  - How will be the weather tomorrow?: “will it be **sunny**, **cloudy** or **rainy**”?  $y \in \{0, 1, 2\}$
- In this case, a best fitting line is not enough
  - Even though this line will be the basis of our **classification model**





# Logistic Regression: Classification

- The obvious idea will be to define a threshold at the classifier output  $h_{\theta}(x)$ , that binarizes the system response:
  - Typically, “0.5” would be the choice, for “*equal classification risks*”
    - It might be more dangerous to predict erroneously one class instead of other one.
    - For example, in a machine learning-based systems for medical diagnosis, classes have different risk.
      - Predict a “malignant cancer” on a “healthy” subject represents a unnecessary concern for the patient and would probably imply to perform additional (an unnecessary) exams.
      - However, provide a “healthy” response for a patient suffering of a “malignant cancer” might represent the patient dead sentence.

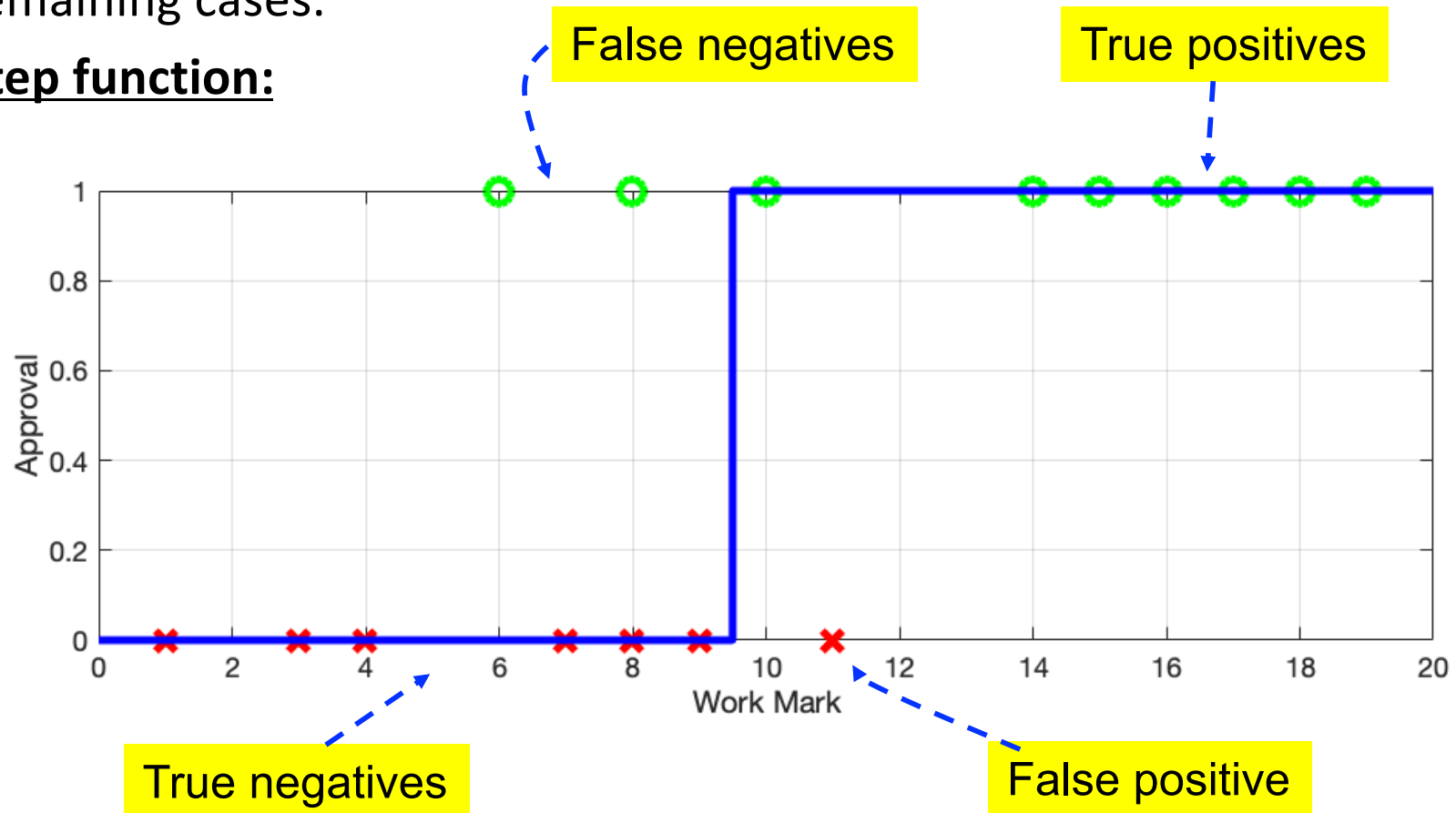
- $$f(x) = \begin{cases} 0, & h_{\theta}(x) < 0 \\ 1, & h_{\theta}(x) \geq 0 \end{cases}$$

- Hence, the response of our classification system can be seen as a composition of two functions:  $f = g \circ h$ 
  - “*f*” is “*g*” after “*h*”
- *We have seen “h” before, but what is “g”?*

# Logistic Regression: Classification

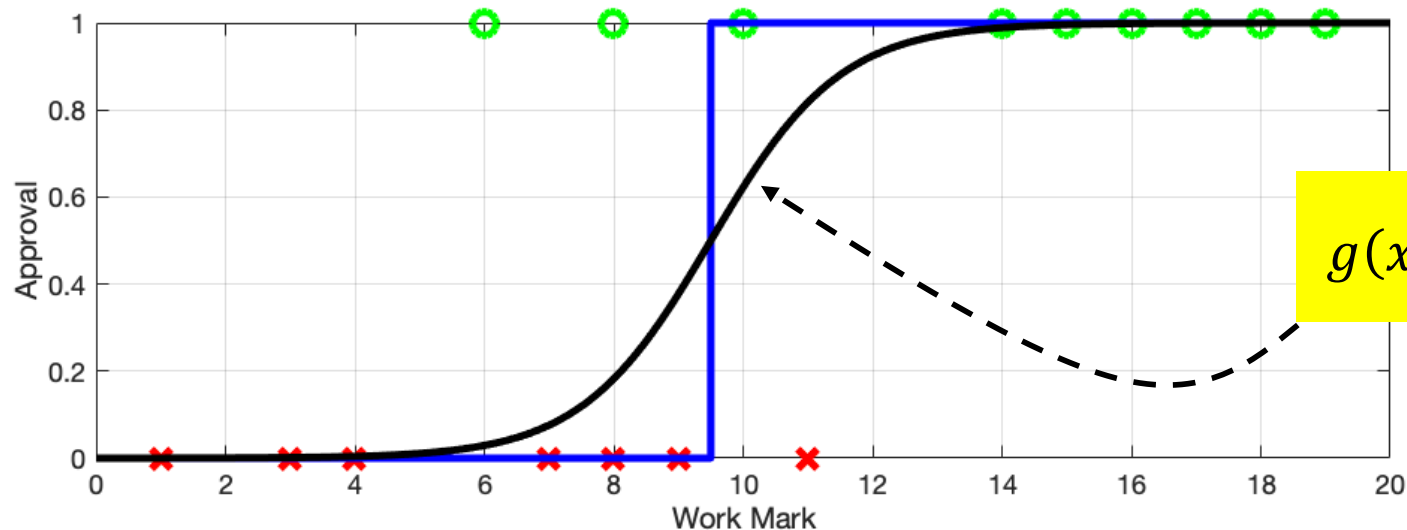
- Essentially, “g” performs a binarization of its input, and produces “1” responses when the input is higher than some threshold, and “0” in the remaining cases.

- Step function:**



# Logistic Regression: Classification

- Assuming the step function as “g”, and  $f = g \circ h$ , obtaining the automatic optimal parameterization of “f” with respect to our data (i.e., machine learning) yields two problems:
  - **Problem 1:** “g” is **not differentiable**
    - It has not a continuous derivative at a single point
  - **Problem 2:** in every other points “g” **has derivative 0**
- The solution is to use a function is close to the step function, without suffering of the above described problems.
  - **Sigmoid Function**



# Logistic Regression: Classification

- Using this composition of functions, our classification system is given by:

$$f_{\theta}(x) = \frac{1}{1 + e^{h_{\theta}(x)}}$$

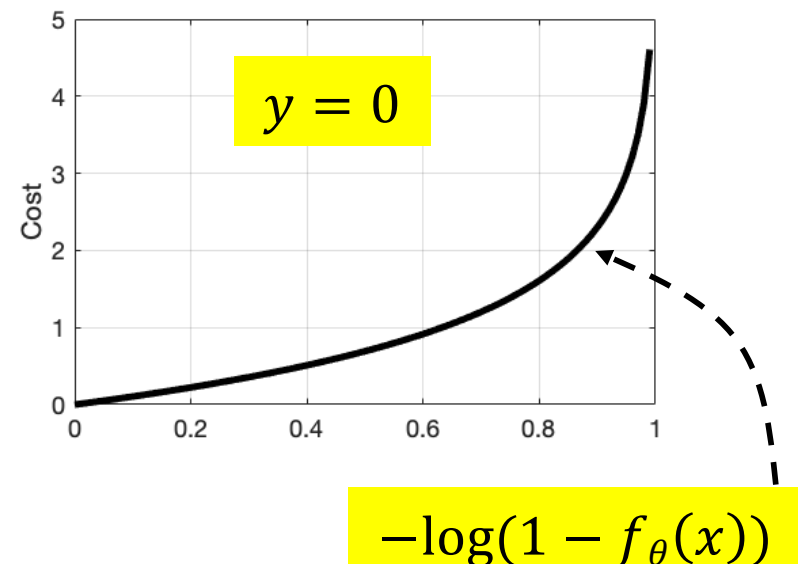
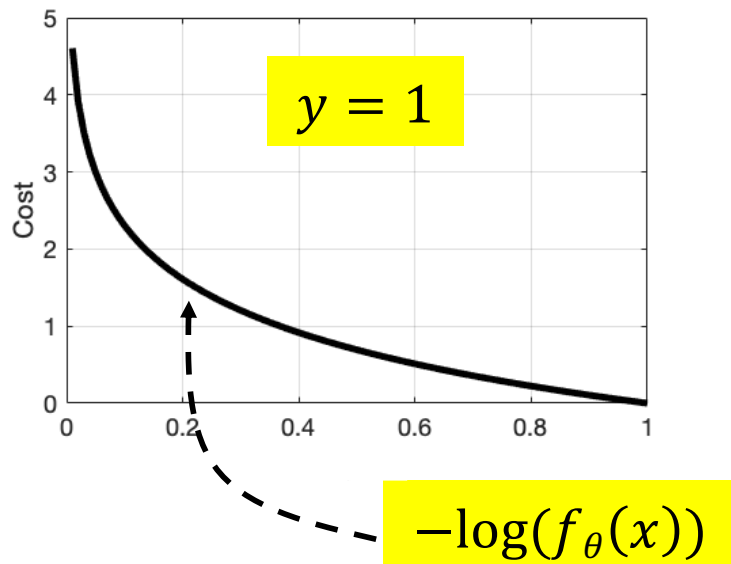
- Or:

$$f_{\theta}(x) = \frac{1}{1 + e^{\theta_1 x + \theta_2}}$$

- The remaining problem is the same as in linear regression:
  - How to find the  $\theta$  optimal parameterization?
- According to the basic principles of Machine Learning, up to now, we've only defined our model.
  - It is also required to define a "Cost Function" (Loss function) that measures how good it is na hypothesis.
  - And a systematic optimization process

# Logistic Regression: Cost Function

- As previously, the cost function will measure how well the model responses ( $f_{\theta}(x)$ ) resemble the “ground-truth” ( $y$ )
  - Intuitively, in cases where the system is supposed to output a “1” and the model predicts a “1”, the cost should be “0”.
  - The same thing should hold for “0” responses.
  - However, the cost (loss) should grow in cases when the system response is far from the ground-truth.
    - The  $\log()$  function is a good choice for representing the desired costs (losses)
    - It varies non-linearly with respect to the distance between the desired and actual responses
      - Attempts to avoid “ridiculously wrong responses”.



# Logistic Regression: Cost Function

- Hence, the cost function for one instance is given by:

$$\bullet \text{ Cost}(f_{\theta}(x), y) = \begin{cases} -\log(f_{\theta}(x)), & y = 1 \\ -\log(1 - f_{\theta}(x)), & y = 0 \end{cases}$$

- And the cost function for the whole dataset is given by the sum of the individual costs:

$$\mathbf{J}(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N (\text{Cost}(f_{\theta}(x^{(i)}), y^{(i)}))$$

- Considering that  $y$  can only assume 2 values (0 or 1), we have:

$$\mathbf{J}(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N y^{(i)} \log(f_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\theta}(x^{(i)}))$$

# Logistic Regression: Optimization

- The optimization can be done exactly as in the linear regression case.
- Using the gradient descent strategy, it is required to find the derivatives of the cost function  $J()$  with respect to the  $\theta$  parameters:

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

- In matrix form, we have:

$$\theta = [\theta_0, \theta_1]^T$$

$$\mathbf{x}^{(i)} = [x^{(i)}, 1]^T$$

- $f_{\theta}(\mathbf{x}) = \frac{1}{1+e^{-\theta^T \mathbf{x}}}$

- $\log(f_{\theta}(\mathbf{x})) = \log\left(\frac{1}{1+e^{-\theta^T \mathbf{x}}}\right)$   
 $= -\log\left(\frac{1+e^{-\theta^T \mathbf{x}}}{1}\right)$

- $\log(1 - f_{\theta}(\mathbf{x})) = -\theta \mathbf{x} - \log\left(\frac{1+e^{-\theta^T \mathbf{x}}}{1}\right)$

# Logistic Regression: Optimization

- Plugging the two simplified expressions in the original cost function, we obtain:

$$J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N -y^{(i)} \log(1 + e^{-\boldsymbol{\theta}x}) + (1 - y^{(i)}) (-\boldsymbol{\theta}x - \log(1 + e^{-\boldsymbol{\theta}x}))$$

- Which can be simplified to:

$$J(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N y^{(i)} \boldsymbol{\theta}x - \log(1 + e^{-\boldsymbol{\theta}x})$$

- Now, as

$$\frac{\partial}{\partial \theta_j} y^{(i)} \boldsymbol{\theta}x = y^{(i)} \boldsymbol{\theta}x$$

$$\frac{\partial}{\partial \theta_j} \log(1 + e^{\boldsymbol{\theta}x}) = \frac{x_j e^{-\boldsymbol{\theta}x}}{1 + e^{\boldsymbol{\theta}x}} = x_j^i f_{\boldsymbol{\theta}}(x)$$

- We have:

$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \sum_{i=1}^N x_j^i (f_{\boldsymbol{\theta}}(x^{(i)}) - y^{(i)})$$



# Logistic Regression: Multi-class

- Up to now, we've only considering binary classification problems.
- When the number of classes ( $c$ ) is higher than 2, the typical approach is to train " $c$ " classifiers
  - In each classifier  $f_{\theta}^{(i)}(x)$ , instances of the  $i^{\text{th}}$  class are considered positive examples, whereas instances of all the remaining classes are treated as negative instances.
- During classification, we pick the class that produces the maximum output response, i.e.:

$$\max_i f_{\theta}^{(i)}(x)$$

