# Fundamentos em Redes Neurais

Regressão Logística

Profº - Dr. Thales Levi Azevedo Valente

thales.l.a.valente@gmail.com.br
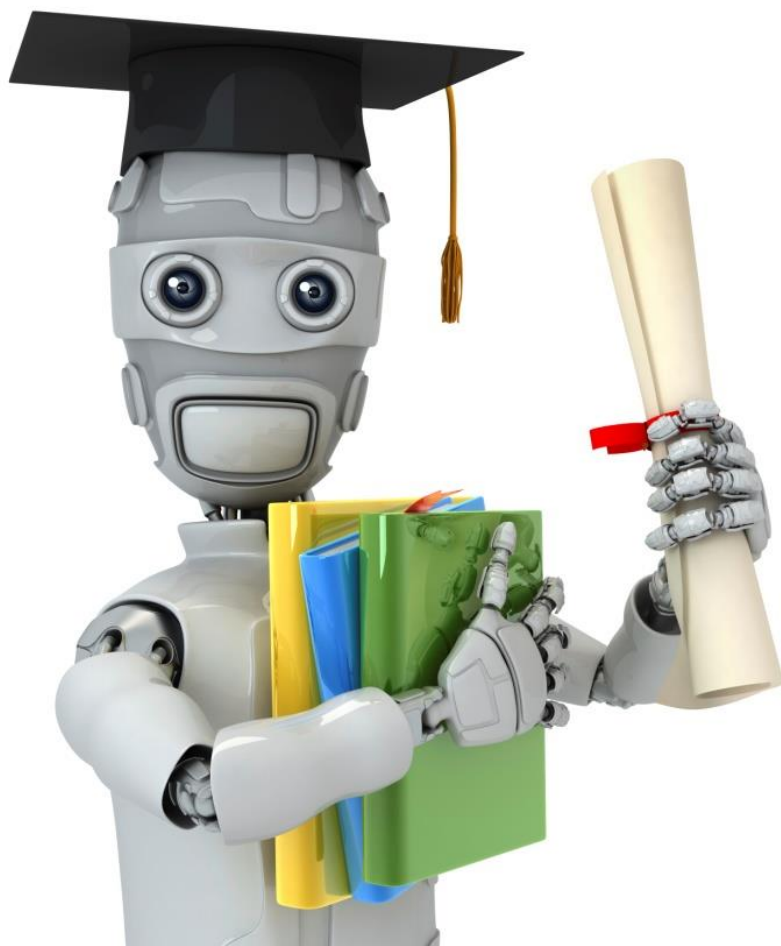
**Governo Federal**

# Sejam Bem-vindos !



**Os celulares devem ficar no silencioso ou desligados**

Pode ser utilizado apenas em caso de emergência



**Boa tarde/noite, por favor e com licença DEVEM ser usados**

Educação é essencial

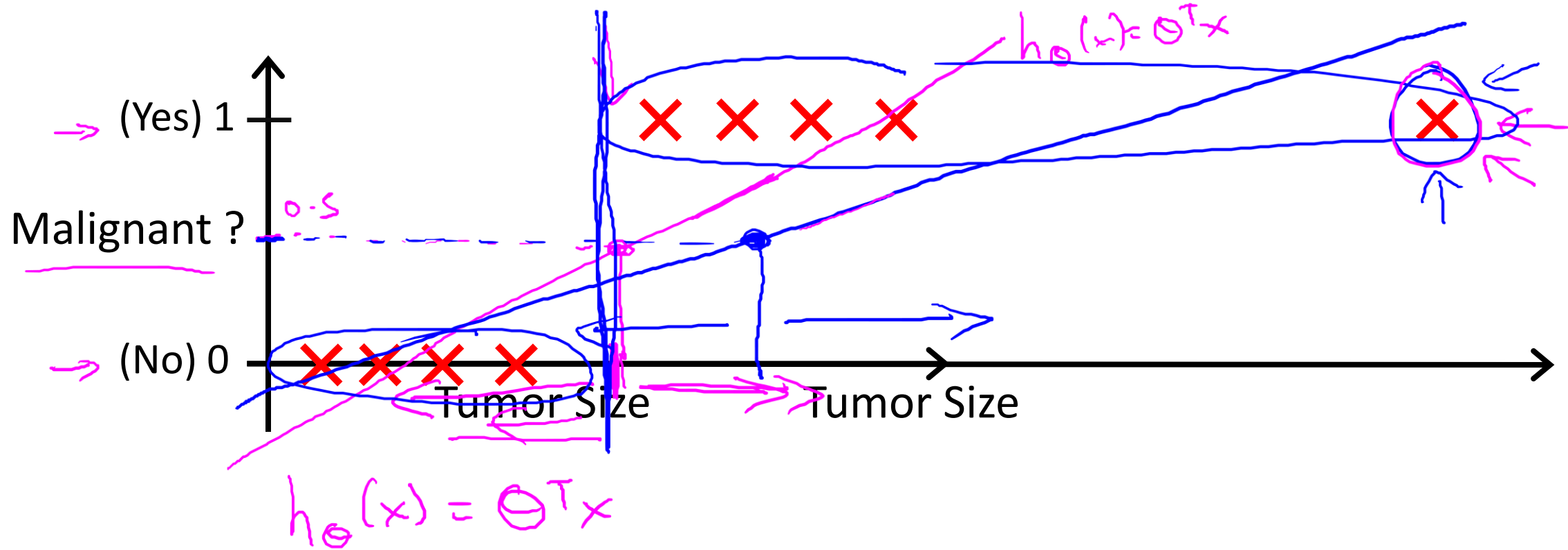# Logistic Regression

## Classification

Machine Learning

# Classification

Email: Spam / Not Spam?

Online Transactions: Fraudulent (Yes / No)?

Tumor: Malignant / Benign ?

$$y \in \{0, 1\}$$

0: "Negative Class" (e.g., benign tumor)

1: "Positive Class" (e.g., malignant tumor)

$$y \in \{0, 1, 2, 3\}$$

Threshold classifier output $h_\theta(x)$ at 0.5:
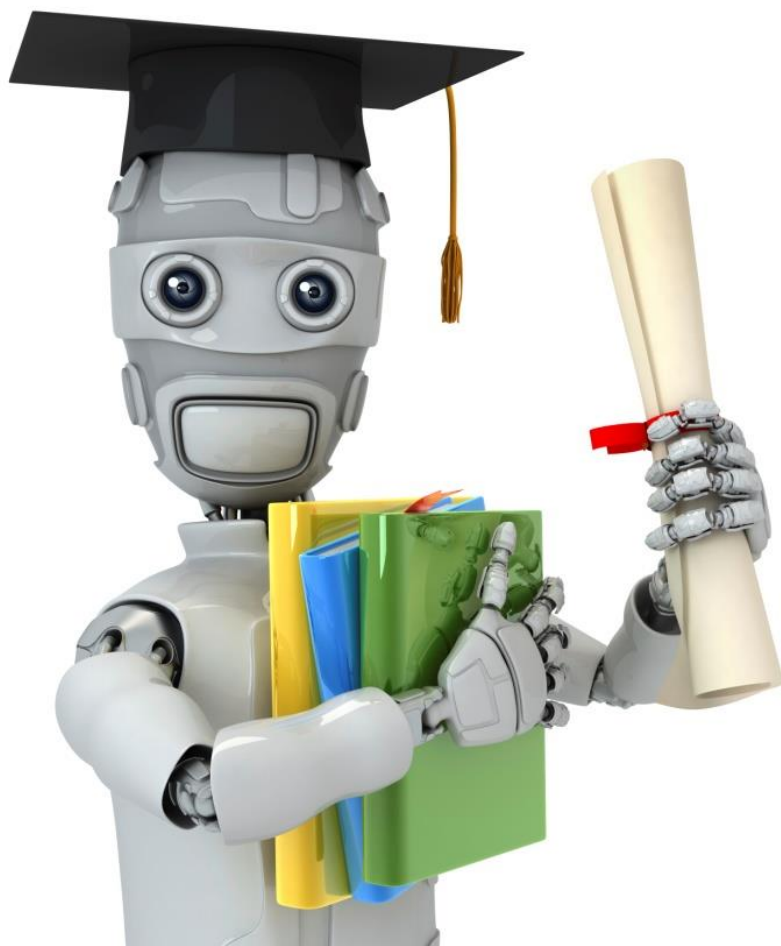
If $h_\theta(x) \geq 0.5$, predict "y = 1"

If $h_\theta(x) < 0.5$, predict "y = 0"

Classification:   y  =  0  or  1

$h_\theta(x)$ can be > 1 or < 0

Logistic Regression:   $0 \le h_\theta(x) \le 1$

Classification

# Logistic Regression

## Hypothesis Representation

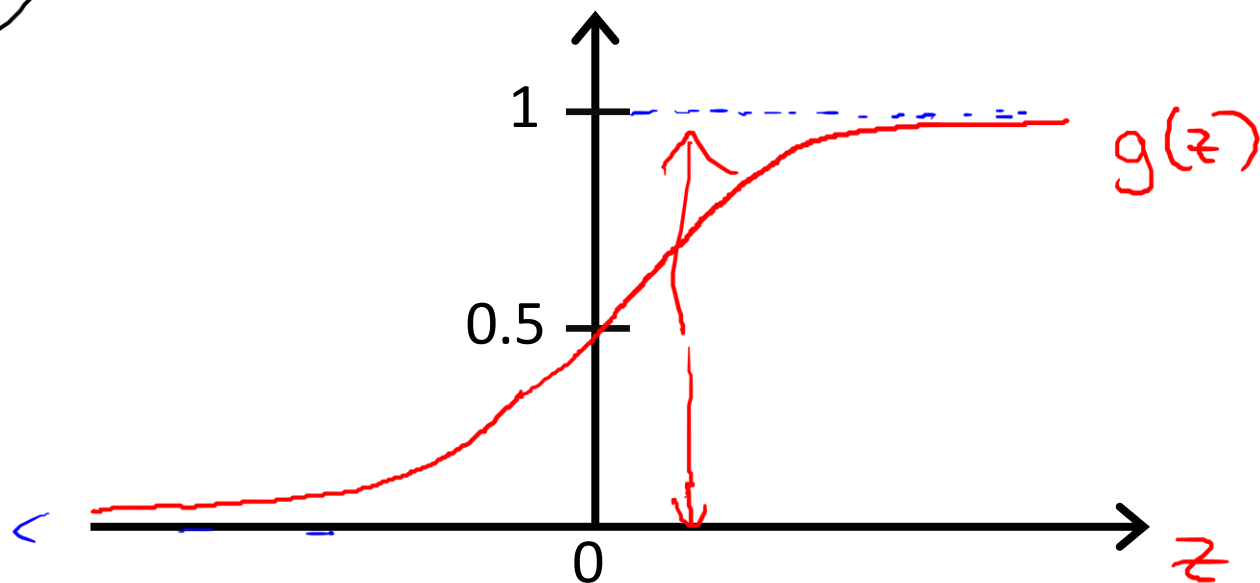Machine Learning

# Logistic Regression Model

Want $0 \le h_\theta(x) \le 1$

$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$\theta^T x$

Sigmoid function
Logistic function

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



$g(z)$

Parameters $\theta$.

# Interpretation of Hypothesis Output

$h_\theta(x)$

$h_\theta(x)$ = estimated probability that $y = 1$ on input x

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$
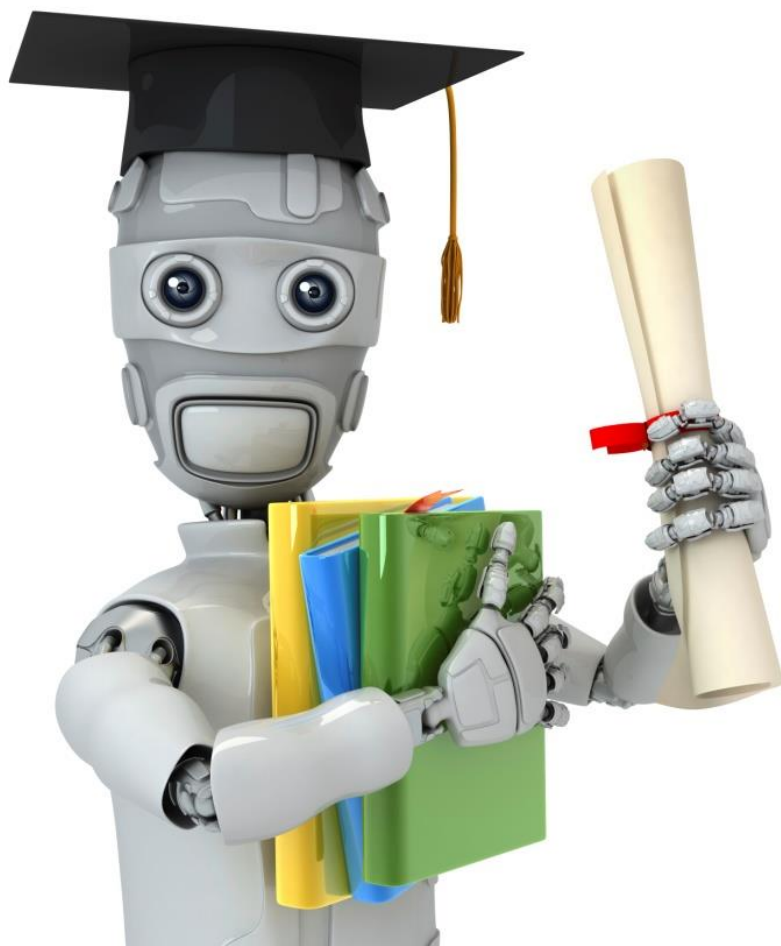
$h_\theta(x) = 0.7$     $y = 1$

Tell patient that 70% chance of tumor being malignant

$h_\theta(x) = P(y=1|x;\theta)$     "probability that y = 1, given x, parameterized by $\theta$"

$y = 0$ or $1$

$\rightarrow P(y=0|x;\theta) + P(y=1|x;\theta) = 1$

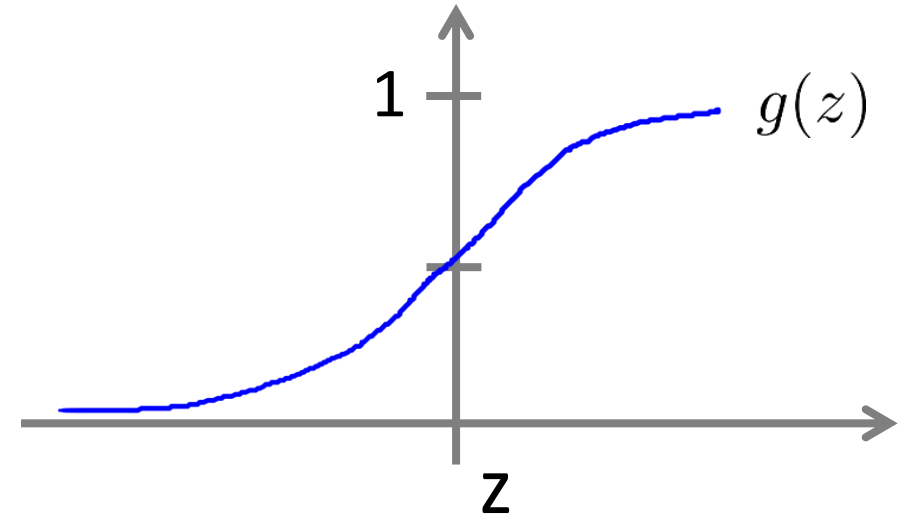$\rightarrow P(y=0|x;\theta) = 1 - P(y=1|x;\theta)$

# Logistic Regression

Decision boundary

Machine Learning

**Logistic regression**

$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$



Suppose predict "$y = 1$" if $h_\theta(x) \geq 0.5$

$\theta^T x \geq 0$

predict "$y =$ "0 if $h_\theta(x) < 0.5$

$\theta^T x < 0$

$g(z) \geq 0.5$
when $z \geq 0$

$h_\theta(x) = g(\theta^T x$

$g(z) < 0.5$
when $z < 0$
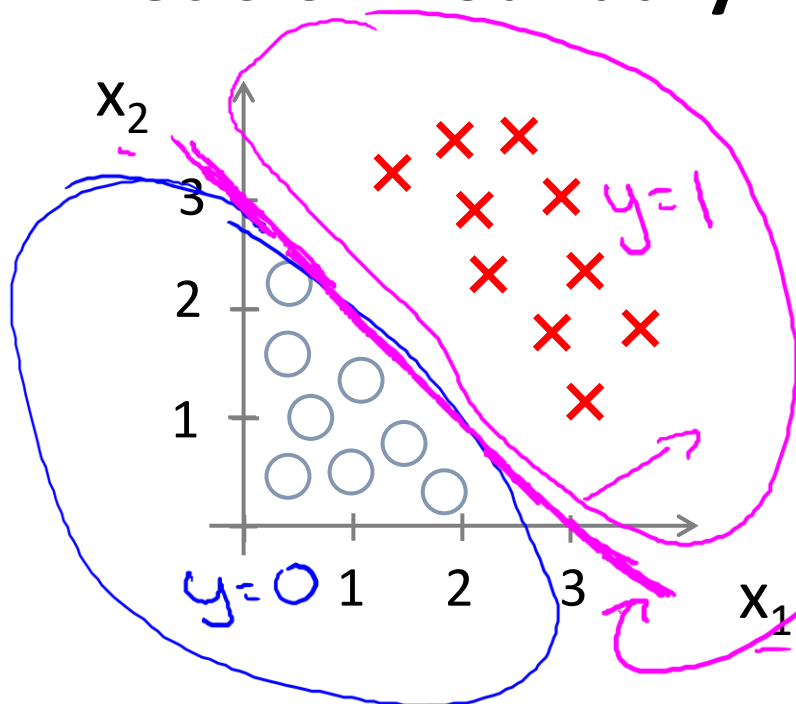
# Decision Boundary



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\underset{-3}{=} \quad \underset{1}{=} \quad \underset{1}{=}$$

Decision boundary

Predict " $y = 1$ " if  $-3 + x_1 + x_2 \geq 0$

$\theta^T x$

$\to x_1 + x_2 \geq 3$

$x_1, x_2$

$\to h_\theta(x) = 0.5$

$x_1 + x_2 = 3$

$x_1 + x_2 < 3$
$\to y = 0$

# Non-linear decision boundaries



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Predict "$y = 1$" if $-1 + x_1^2 + x_2^2 \geq 0$

$$x_1^2 + x_2^2 = 1$$

$$x_1^2 + x_2^2 \geq 1$$

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \ldots)$$

Logistic
Regression

Cost function

Machine Learning

Training set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \cdots, (x^{(m)}, y^{(m)})\}$

m examples $\qquad x \in \begin{bmatrix} x_0 \\ x_1 \\ \cdots \\ x_n \end{bmatrix}$ $\mathbb{R}^{n+1}$ $\qquad x_0 = 1, y \in \{0, 1\}$

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

How to choose parameters $\theta$ ?

# Cost function

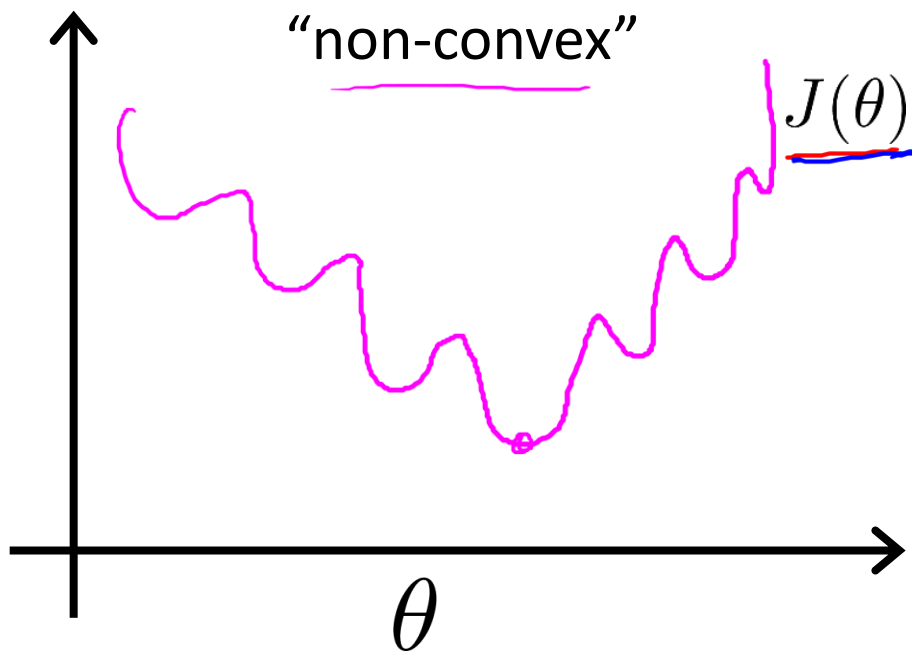Linear regression: $\quad J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$

Logistic

$\text{Cost}(h_\theta(x^{(i)}), y)$

$$\text{Cost}(h_\theta(x^{(i)}), y^{(i)}) = \frac{1}{2} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$$\frac{1}{1 + e^{-\theta^T x}}$$

"non-convex"

$J(\theta)$

"convex"

$J(\theta)$

$\theta$

$\theta$

# Logistic regression cost function

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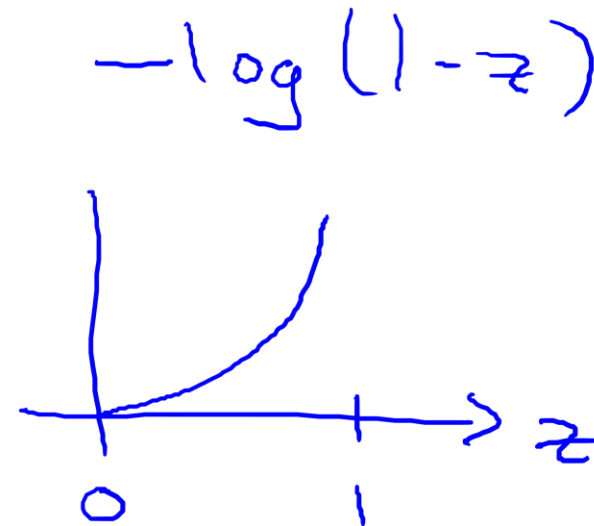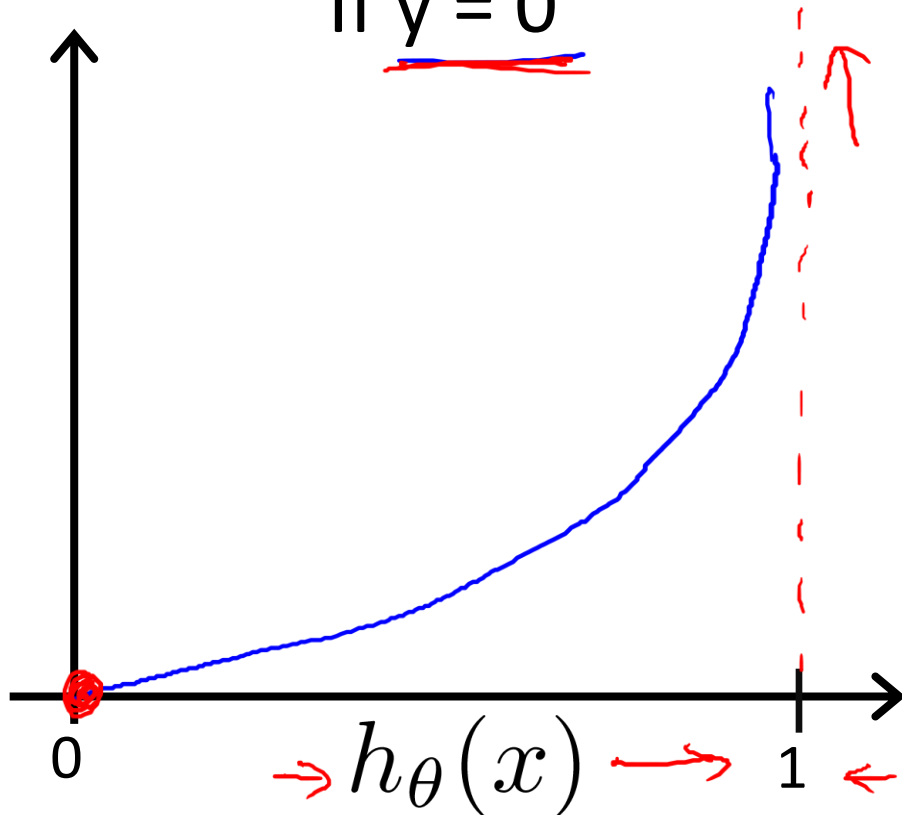
If y = 1

$\text{Cost} = 0$ if $y = 1, h_\theta(x) = 1$

But as $\quad h_\theta(x) \to 0$

$Cost \to \infty$

Captures intuition that if $h_\theta(x) = 0$, (predict $P(y = 1 | x; \theta) = 0$), but $y = 1$, we'll penalize learning algorithm by a very large cost.

$h_\theta(x)$

0          1

# Logistic regression cost function

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

If y = 0



$-\log(1-z)$

$0 \rightarrow h_\theta(x) \rightarrow 1 \leftarrow$

# Logistic Regression

Simplified cost function and gradient descent

Machine Learning

# Logistic regression cost function

$$\rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\rightarrow \text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \leftarrow \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \leftarrow \end{cases}$$

Note: $y = 0$ or $1$ always

$$\rightarrow \text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1 - h_\theta(x)) \leftarrow$$

$$= 0 \qquad = 0 \qquad = 1$$

If $y = 1$: $\text{Cost}(h_\theta(x), y) = -\log h_\theta(x) \leftarrow$

If $y = 0$: $\text{Cost}(h_\theta(x), y) = -\log(1 - h_\theta(x))$

# Logistic regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} [\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))]$$

To fit parameters $\theta$:

$$\min_\theta J(\theta) \qquad \text{Get } \theta$$

To make a prediction given new $x$:

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

$$p(y=1 \mid x; \theta)$$

**Gradient Descent**

$$\rightsquigarrow J(\theta) = -\frac{1}{m}[\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)}))]$$

Want $\min_\theta J(\theta)$:

Repeat {

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

(simultaneously update all $\theta_j$)

}

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

# Gradient Descent

$$J(\theta) = -\frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)}))\right]$$

Want $\min_\theta J(\theta)$:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \begin{matrix} \leftarrow \\ \leftarrow \\ \leftarrow \\ \\ \leftarrow \end{matrix}$$

for $i = 0$ to $n$

Repeat {

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

(simultaneously update all $\theta_j$)

}

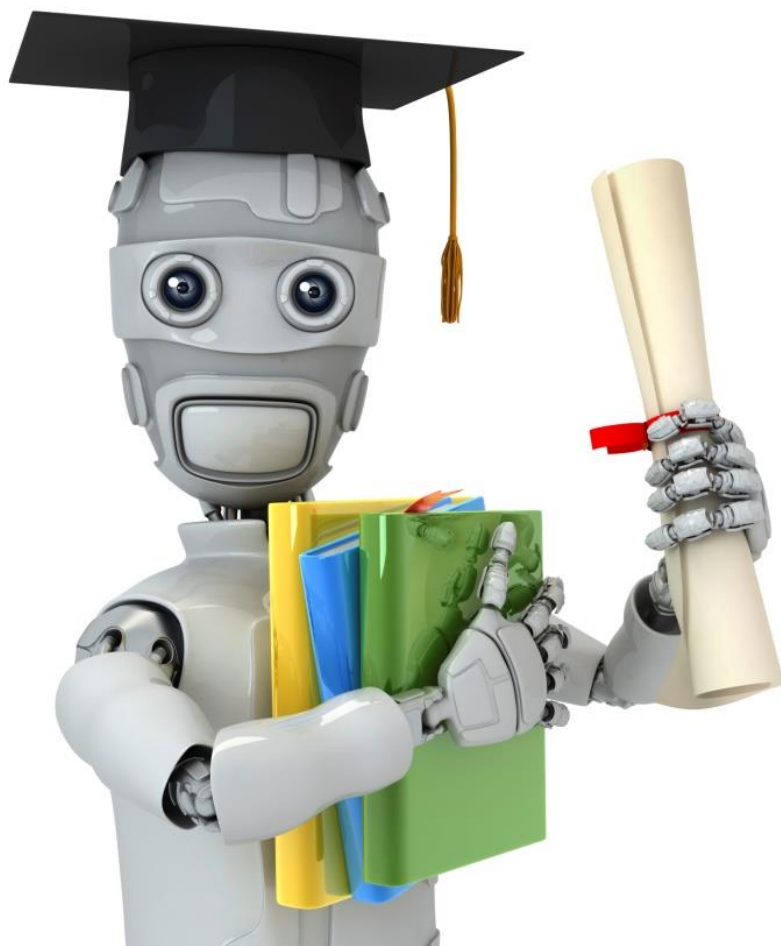$h_\theta(x) = \theta^T x$

$h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$

Algorithm looks identical to linear regression!

# Logistic Regression

## Advanced optimization

Machine Learning

# Optimization algorithm

Cost function $J(\theta)$. Want $\min_\theta J(\theta)$.

Given $\theta$, we have code that can compute

- $J(\theta)$
- $\frac{\partial}{\partial \theta_j} J(\theta)$      (for $j = 0, 1, \ldots, n$ )

Gradient descent:

Repeat $\{$

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

$\}$

**Optimization algorithm**

Given $\theta$, we have code that can compute

- $J(\theta)$
- $\frac{\partial}{\partial \theta_j} J(\theta)$      (for $j = 0, 1, \ldots, n$)

Optimization algorithms:
- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:
- No need to manually pick $\alpha$
- Often faster than gradient descent.

Disadvantages:
- More complex

$$\text{theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \begin{matrix} \text{--- theta(1)} \leftarrow \\ \text{--- theta(2)} \\ \\ \text{--- theta(n+1)} \end{matrix}$$

```
function [jVal, gradient] = costFunction(theta)

  jVal = [code to compute J(θ)];

  gradient(1) = [code to compute  ∂/∂θ₀ J(θ)];

  gradient(2) = [code to compute  ∂/∂θ₁ J(θ)];

       .
       .
       .

  gradient(n+1) = [code to compute  ∂/∂θₙ J(θ)   ];
```

$\text{jVal} = [\text{code to compute } J(\theta)];$

$\text{gradient(1)} = [\text{code to compute } \frac{\partial}{\partial \theta_0} J(\theta)];$

$\text{gradient(2)} = [\text{code to compute } \frac{\partial}{\partial \theta_1} J(\theta)];$

$\text{gradient(n+1)} = [\text{code to compute } \frac{\partial}{\partial \theta_n} J(\theta)];$

# Logistic Regression

Multi-class classification: One-vs-all

Machine Learning

# Multiclass classification

Email foldering/tagging: Work, Friends, Family, Hobby

$y=1 \quad y=2 \quad y=3 \quad y=4$
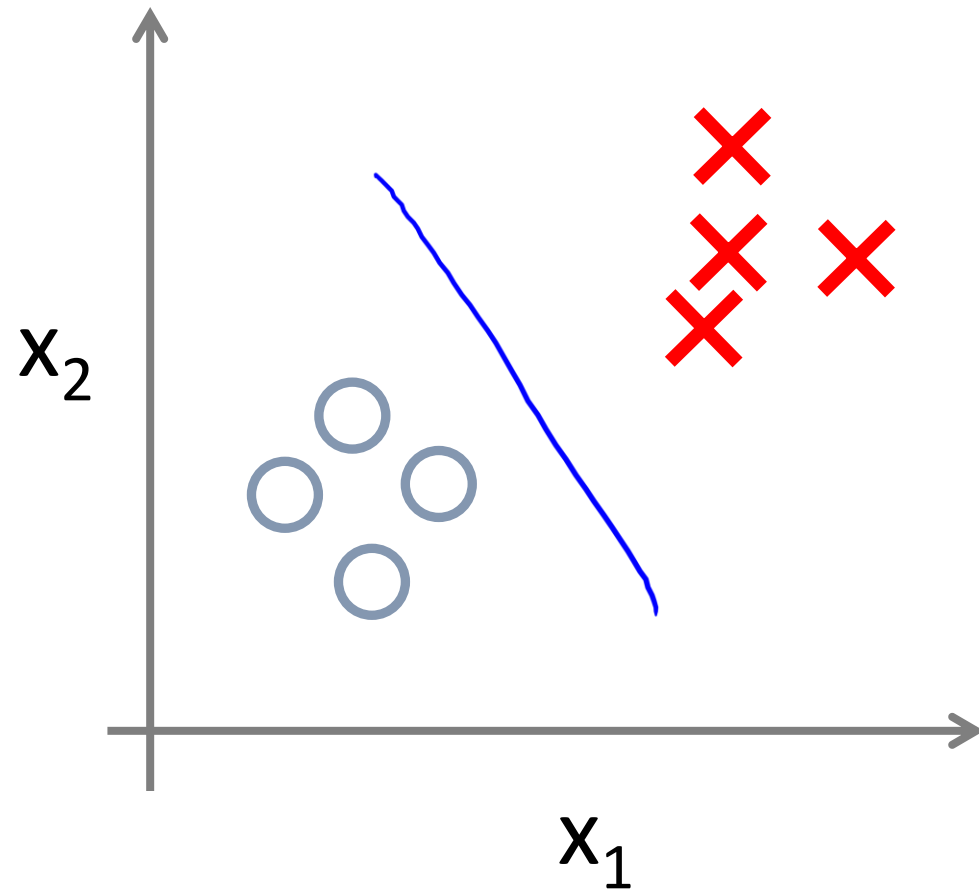
Medical diagrams: Not ill, Cold, Flu

$y=1 \quad 2 \quad 3$

Weather: Sunny, Cloudy, Rain, Snow

$y=1 \quad 2 \quad 3 \quad 4 \leftarrow$

0 — 1 — 2 — 3

Binary classification:

Multi-class classification:

# One-vs-all (one-vs-rest):



Class 1: △ ←

Class 2: □ ←

Class 3: ✕ ←

$$h_\theta^{(i)}(x) = P(y = i | x; \theta) \qquad (i = 1, 2, 3)$$

## One-vs-all

Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class $i$ to predict the probability that $y = i$.

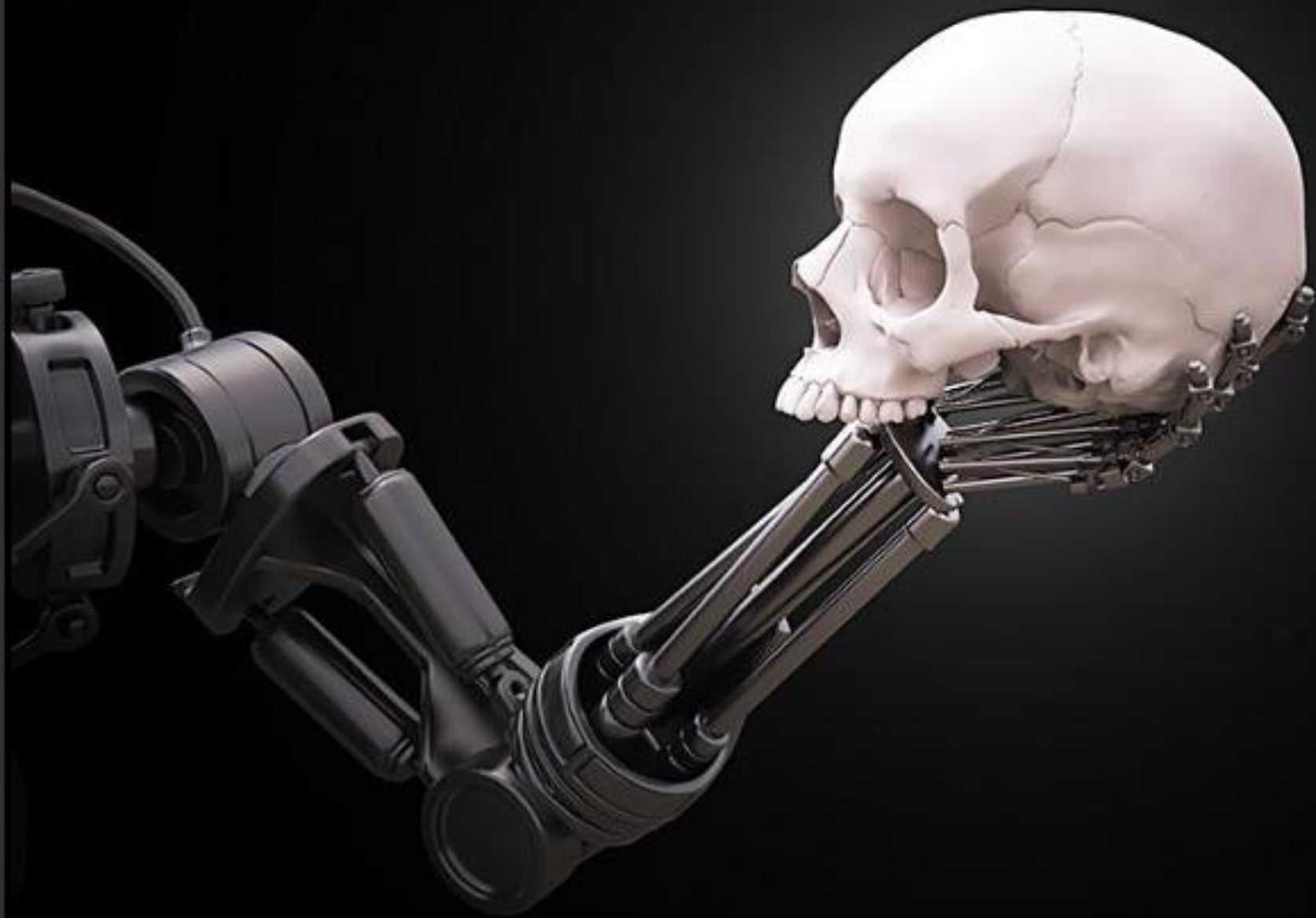On a new input $x$, to make a prediction, pick the class $i$ that maximizes

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

## Referências

NG, Andrew. *Machine Learning*. Stanford University, 2011. Curso oferecido via Coursera. Disponível em: https://www.coursera.org/learn/machine-learning.

Dúvidas?

Até a próxima...

Apresentador

# Thales Levi Azevedo Valente

E-mail:
thales.l.a.valente@gmail.com