



**MÓDULO 2**

**Aprendizagem supervisionada**

Especialização em Machine Learning

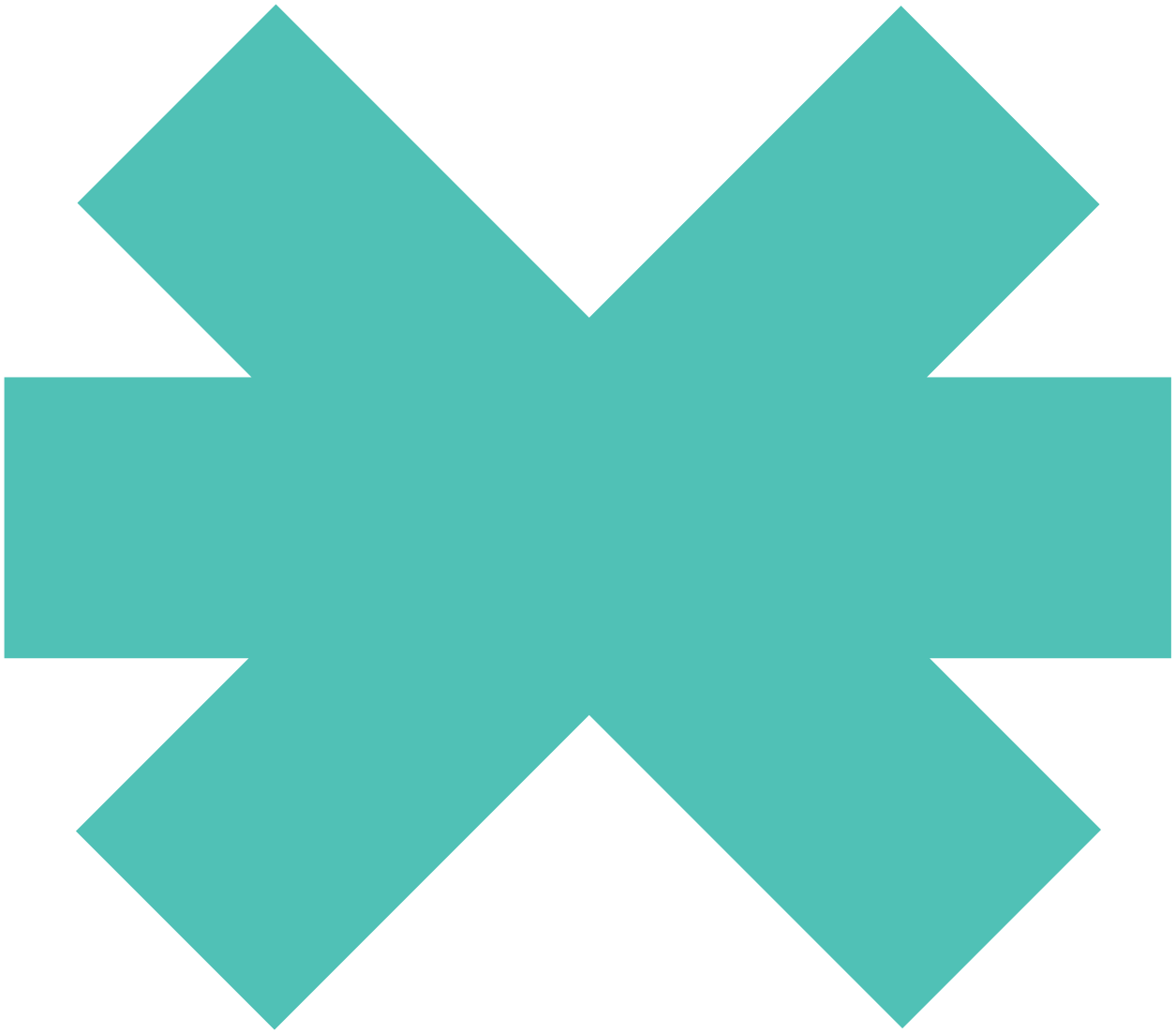
1

# Regressão linear



New  
Technology  
School

**Tokio.**



# 1 Regressão linear

## Sumário

1.1	Regressão linear simples	05
1.2	Regressão linear múltipla	08

A aprendizagem automática consiste na modelação de sistemas, através do treino de um modelo que aprende com um conjunto de dados, cujos resultados são previamente conhecidos.

Depois desse treino inicial, obtemos o modelo, os pesos, os coeficientes, etc., que poderemos usar para prever ou descrever novos exemplos, cujos resultados ainda desconhecemos.

Assim, trata-se de treinar um modelo com perguntas, conhecendo as suas respostas para, posteriormente, poder usá-lo para resolver perguntas cujas respostas desconhecemos.

### APLICAÇÕES

- Preditores ou avaliadores de preços de venda de qualquer produto, como imóveis, automóveis, bens, etc.
- Preditores de procura, oferta, níveis de produção, etc.
- Classificação de frutos segundo classes, qualidades, gamas e usos.
- Deteção de *spam* em *emails*, comentários, contas *fake*, etc.
- Manutenção preditiva, predição de falhas.
- Visão computacional.
- Reconhecimento de voz, agentes virtuais, transcrição de texto por voz.
- Processamento da linguagem natural.

## 1.1 Regressão linear simples

Procuramos uma forma de modelar um conjunto de dados, de propor uma função que relacione linearmente uma variável independente (os dados dos exemplos) com uma variável dependente (os seus resultados, a variável a predizer).

Essa função modela o conjunto de dados, isto é, simula o seu comportamento de uma forma matemática, que ajude a explicá-lo e a realizar previsões no futuro.

### Equação do modelo

A equação linear simples é a equação da reta:

$$Y = h_{\Theta}(x) = \Theta^T \times x = \theta_0 + \theta_1 \times x_1$$
$$x_0 = 1$$

Figura 1

Código Latex:

```
Y=h_{\Theta}(x)=\Theta^T \times x=\theta_0+\theta_1 \times x_1 \\ x_0=1
```

A equação representa a hipótese que existe de uma relação linear, entre duas variáveis,  $X$  e  $Y$ , modelada desta forma e dependente de  $\Theta$ .

- $x$  poderia representar o nível socioeconômico de um bairro.
- $Y$  poderia representar o nível de delinquência desse bairro.
- $\theta_0$  seria o “parâmetro de intersecção” ou nível base do modelo.
- $\theta_1$  seria o coeficiente da variável  $x$  que explica como varia o nível de delinquência, em função do nível socioeconômico do bairro.
- $x_0$ , por conveniência, será 1; desta forma, podemos multiplicar diretamente ambos os vetores, em vez de multiplicar um vetor e realizar uma soma.

## Representação gráfica

Y	x1
-12,42	-2
-7,35	-1
10,64	0
5,43	1
15,35	2
23,53	3
50,13	4
64,05	5
81,21	6
81,03	7
101,51	8
116,16	9
87,34	10
99,00	11
99,38	12
136,59	13
127,14	14
161,28	15
173,39	16
192,36	17
174,91	18

Y vs x1

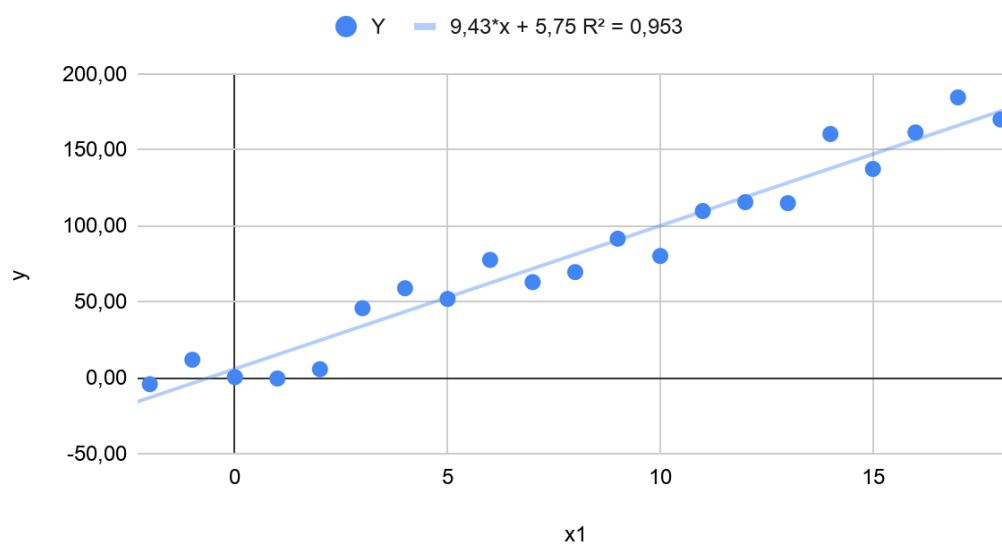


Figura 2

## Tipos de variáveis

As variáveis podem ser de vários tipos:

- **Quantitativas.** Geralmente trata-se de variáveis de tipo numérico que representam uma característica mensurável.
  - **Contínuas:** variáveis numéricas em números reais, isto é, que podem apresentar decimais. Por exemplo: preços, temperatura, peso, distâncias, etc.
  - **Cardinais:** variáveis numéricas em números inteiros, isto é, que não podem apresentar decimais e que se representam com *int*, *double*, etc. Por exemplo: n.º de produtos, de familiares, de compras realizadas, etc.
- **Qualitativas ou categóricas.** Variáveis que não são de tipo numérico, mas representam uma propriedade, etiqueta ou categoria atribuída ao exemplo.
  - **Ordinais:** categorias que podem e devem ser ordenáveis, uma vez que existe uma certa disposição ou hierarquia natural que devemos obter. Por exemplo: o nível socioeconómico ("baixo, médio, alto"), nível educativo ("primária, secundária, universidade, mestrado, doutoramento"), hierarquia ("júnior, sénior, agente, diretor", "soldado, capitão, tenente, general") etc.
  - **Nominais:** categorias não ordenáveis. Por exemplo: cores, género, marcas, os ID de produto, "louro, moreno, ruivo", etc.
  - **Binárias:** categorias não ordenáveis ou nominais que dispõem apenas de dois valores. Por exemplo: "sim/não", "homem/mulher", "cliente existente/não", "utilizador de pagamento/gratuito", "contratou o serviço/não", etc.

## 1.2 Regressão linear múltipla

Uma **regressão linear múltipla** constitui, tal como a regressão linear simples, um modelo estatístico que estima uma previsão numérica de uma variável objetivo ou dependente, com base na combinação linear de variáveis.

Na regressão linear simples apenas temos em conta uma única variável, enquanto numa múltipla temos em consideração mais do que uma, para verificar como essa variável afeta o objetivo.

### Aplicações:

- Predição de preços de venda.
- Previsão de demanda ou nº de vendas.
- Predição de probabilidade de compra.
- Manutenção preditiva de maquinaria: predição de probabilidade de falhas, tempo até à próxima falha, nível de desgaste, etc.
- Previsões meteorológicas.

### Modelação de dados

O nosso objetivo é recolher, para cada exemplo que usarmos durante o treino, todas as características que criarmos, que possam influenciar a variável objetivo e dependente ( $Y$ ) e codificá-las na variável  $X$ .

A variável  $X$  contará, assim, com uma linha, por exemplo, com um valor por cada uma das características recolhidas nesse exemplo.

#### EXEMPLO: PERITO AVALIADOR DE IMÓVEIS

Para tentar treinar um modelo de ML, face a prever o preço de venda de um imóvel, devemos planejar primeiro, qual será a nossa variável objetivo.

Para este caso, será o preço de venda, mas em €? Em milhares de €, com cêntimos de € ou sem decimais? Numa escala distinta?

Estas são algumas das perguntas que podemos fazer, sobre como codificar a variável  $Y$ .

Da mesma forma, seguidamente, iremos planejar os exemplos a considerar: estes serão imóveis cujas características e preços de venda conhecemos (procedentes, por exemplo, da base de dados de uma grande imobiliária), que formarão a variável matriz  $X$ .

Sobre a variável  $Y$  podemos alertar para: o que deveríamos usar? Os preços a que são anunciados os imóveis? Os resultados de um avaliador oficial? Apenas os preços de venda final dos imóveis? Podemos aceder-lhes?



Sobre esses exemplos, devemos considerar que características ou atributos podem afetar o preço final do imóvel e quais deles podemos obter com veracidade e precisão suficientes.

Para este exemplo, podemos planejar as seguintes características e considerações:

$m \times n$	$x_{...0}$	Preço $x_{...1}$	$m^2$ $x_{...2}$	Hab. $x_{...3}$	WC $x_{...4}$	Tipo $x_{...5}$	Garagem $x_{...6}$	Ano $x_{...7}$	Dist. $x_{...8}$	Serv. $x_{...9}$	Loc. $x_{...10}$	Terr. $x_{...11}$	Pisc. $x_{...12}$
Casa 1 $x_{0...}$	1	150	110	3	2	1 (pisso)	1	2010	1274	8,42	1	0	0
Casa2 $x_{1...}$	1	200	100	4	2	3 (ático)	1	2007	3265	6,76	3	2	0
Casa 3 $x_{2...}$	1	175	130	4	2	2 (casa)	2	1988	6312	4,82	5	1	0
Casa 4 $x_{3...}$	1	110	90	2	1	1 (pisso)	0	1977	2252	5,77	2	0	1

Sendo  $m$  o n.º de exemplos ou imóveis que iremos considerar, para treinar, e  $n$  o n.º de características obtidas de cada exemplo.

- Superfície útil, em  $m^2$ : superfície útil construída ou ambas? Em  $m^2$ , em intervalos de  $m^2$ , em “pequeno, médio, amplo”? Em superfície por piso?
- Número de divisões: apenas o n.º de quartos ou de divisões no total? Todas as divisões devem contar da mesma forma ou, por exemplo, o quarto de casal como dupla? E o sótão, ou a lavanderia se existirem?
- Número de casas de banho: contamos com e sem banheira, de igual forma, de forma independente ou a sua superfície?
- Tipo de imóvel: como o representamos? Com *one-hot encoding* (múltiplas características/colunas com valor 0 e apenas um 1 na correta) ou como uma categoria ordenada?
- Garagem: categoria binária ou n.º de lugares de garagem? Superfície? Tipo: “fechado, aberto, individual, comunitário”, etc.?
- Ano: ano de construção ou idade? Anos desde a última remodelação?
- Distância do centro: trata-se de um parâmetro cujo cálculo é muito complexo. Distância linear? Tempo? Apenas em direção ao centro ou a vários destinos habituais, tais como a praia, a saída da autoestrada, centros comerciais e de negócios, etc.?
- Serviços: outro parâmetro difícil de calcular. Podemos propor um valor numérico estimado, com base em vários critérios ou usar esses critérios base, se bem que complicaríamos bastante o modelo e, às vezes, podemos não obter todos os dados.
- Localização: devemos definir o bairro ou área onde se encontra ou um nível comparativo entre várias localizações semelhantes?
- Terraço: categoria binária, n.º de terraços, superfície total? E as varandas?

- Piscina: binário, superfície? Privada, partilhada? Zonas comuns da urbanização?

Assim, neste ponto, devemos tentar definir todas as características de um imóvel, que determinariam quanto um comprador estaria disposto a pagar por ele, o que é bastante complicado e constitui uma fonte de possíveis erros e imprecisões.

O recomendável poderia ser começar com um modelo simples para tentar explorar os seus dados e resultados e, em função desses resultados, continuar a recolher mais informação, com o objetivo de criar um modelo mais preciso e robusto, em cada versão.

## Modelação de uma curva

Em várias ocasiões, necessitamos de modelar conjuntos de dados complexos com relações diferentes, entre as diversas variáveis independentes e a variável objetivo. O que fazemos se não seguirem uma relação linear simples, mas exponencial, potencial, logarítmica, etc.? Podemos resolver esses modelos com uma regressão linear múltipla comum?

Nestes casos, devemos pré-processar as variáveis. Assim, por exemplo, se queremos considerar em polinómio de grau 3, podemos converter a variável original em três variáveis do seguinte modo:

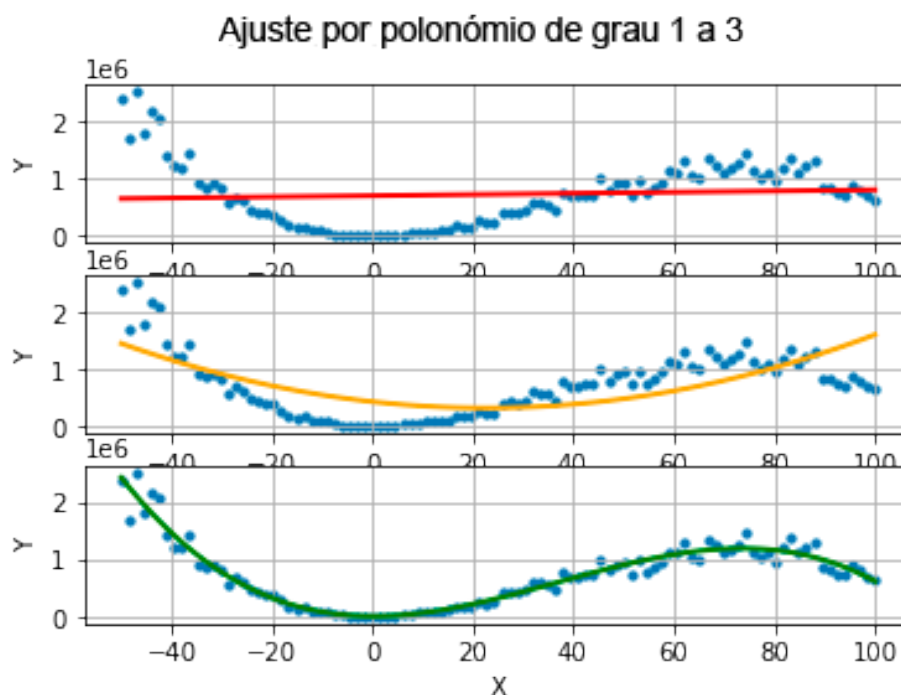


Figura 3

Podemos proceder do mesmo modo para outros tipos de ajustes, como o potencial ou logarítmico:

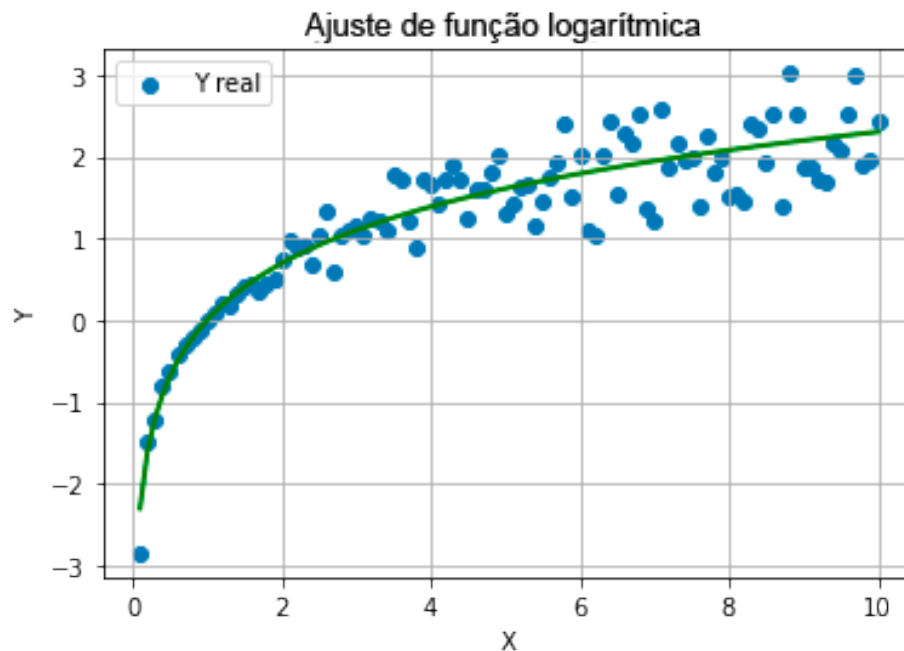


Figura 4

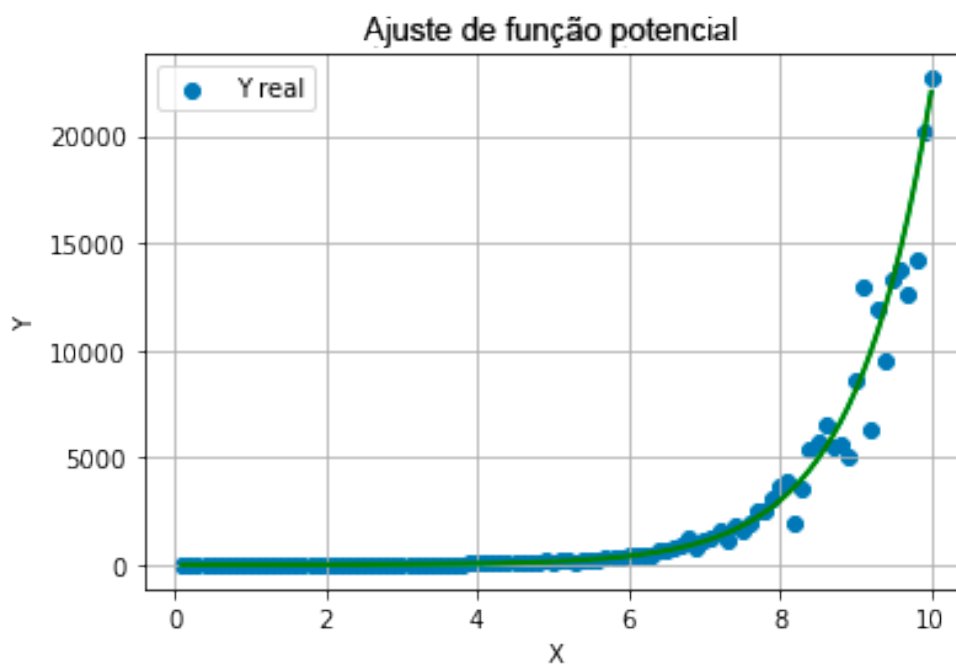


Figura 5

Desta forma, a partir de uma única variável original, consideramos várias variáveis, aparentemente independentes entre si, que são fruto do processamento dessa variável original.

Do mesmo modo, devemos recordar sempre que, segundo as propriedades de um sistema de equações linear, qualquer nova variável que seja uma combinação linear de outras variáveis, não fornecerá nova informação ao sistema ou modelo, logo deverá ser descartada ou simplificada.

## Resolução analítica

Estamos a planear uma relação linear entre as variáveis  $Y$  e  $X$  que, por serem matrizes, nos permitem planear um sistema de equações lineares, que podemos resolver de forma analítica, tal como em qualquer outro problema de álgebra lineal.

Para resolver um sistema de equações lineares, resolvemos a equação seguinte:

$$Y = h_{\Theta}(x) = \Theta^T \times x$$

$$\Theta = (X^T \times X)^{-1} \times X^T \times Y$$

Figura 6

Código Latex:

```
Y=h_{\Theta}(x)=\Theta^T \times x \\
\Theta = (X^T \times X)^{-1} \times X^T \times Y
```

Contudo, surgem, com frequência, vários problemas que nos complicam a sua resolução ou a tornam impossível:

- A matriz resultante não invertível.
- O número de exemplos e características determina que a resolução, por este método, não seja viável por memória ou tempo de execução.
- Recordemos que as equações não são plenamente corretas:  $Y \approx \Theta \times X$ 
  - Estamos a tentar aproximar um modelo, mas pode não seguir realmente essa relação linear, podem faltar características, etc.
  - Os nossos dados podem não estar corretos a 100%.
  - Teremos sempre uma margem de erro ou aleatoriedade.

## Função de custo

Modela o “custo” ou o erro acumulado, da hipótese proposta, para modelar este conjunto de dados sobre todos os exemplos.

O erro será a diferença entre o valor real de  $Y$ , para os exemplos do conjunto de aprendizagem, e a predição de  $Y$  para esses exemplos.

$$Y = h_{\Theta}(x) = \Theta^T \times x$$

$$J_{\theta} = \frac{1}{2m} \sum_{i=0}^m (h_{\theta}(x^i) - y^i)^2$$

Figura 7

Código Latex:

```
Y=h_{\Theta}(x)=\Theta^T \times x \\ J_{\theta} = \frac{1}{2m} \sum_{i=0}^m (h_{\theta}(x^i)-y^i)^2
```

Parâmetros:

- $m$ : n.º de exemplos.
- $i$ : índice do exemplo.
- $x$ : vetor ou linha do exemplo.
- $y$ : resultado conhecido para esse exemplo.

Assim, o objetivo de treinar o modelo é o de obter o valor de  $\theta$  que minimize  $J$ , ou seja, que tenha o menor custo ou erro possível.

---

## Resolução por métodos iterativos

Em vez de uma solução analítica, podemos resolver um sistema de equações lineares através de uma solução iterativa que, ao longo de inúmeras iterações, vai ficando cada vez mais próxima do valor ótimo.

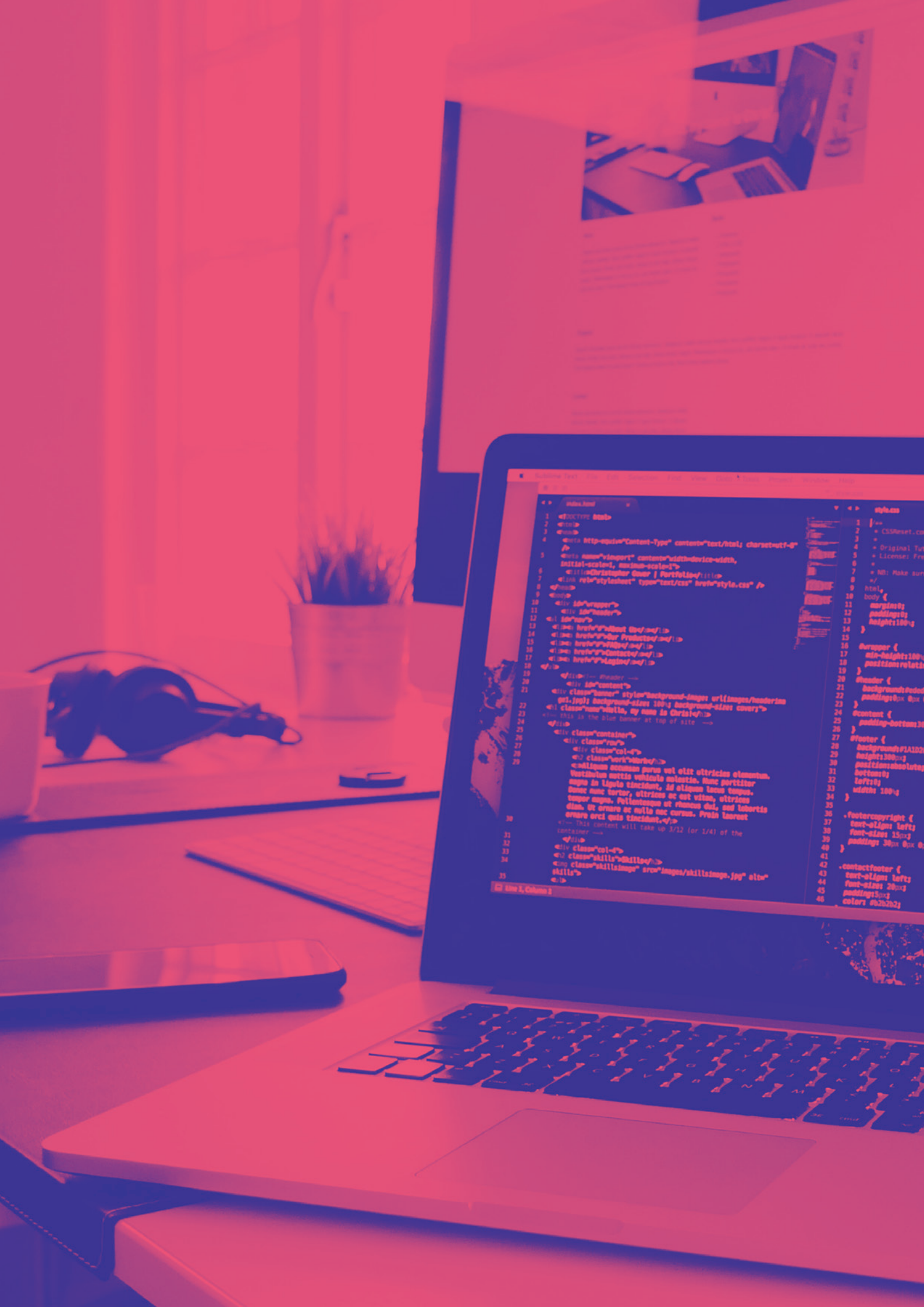
Assim, começemos por propor os valores iniciais de  $\theta$  aleatórios, calculamos a  $J$  para a  $\theta$  e continuamos a atualizar esses valores, aproximando-nos do valor ótimo, ao longo de múltiplas iterações.

Para isso, podemos utilizar qualquer algoritmo de otimização matemática.

---

## Algoritmo de resolução

1. Compilar os exemplos,  $X$ .
2. Compilar os seus resultados conhecidos,  $Y$ .
3. Iniciar aleatoriamente os coeficientes/pesos de  $\theta$ .
4. Iterativamente, calcular o custo  $J$  para cada  $\theta$  e atualizar os seus valores.
5. Finalizar quando  $\theta$  converja num valor final.



```
index.html
1 <!DOCTYPE html>
2 <html>
3 <head>
4 <meta http-equiv="Content-Type" content="text/html; charset=utf-8" />
5 <meta name="viewport" content="width=device-width, initial-scale=1" />
6 <title>Christopher Gnar | Portfolio</title>
7 <link rel="stylesheet" type="text/css" href="style.css" />
8 </head>
9 <body>
10 <div id="wrapper">
11 <div id="header">
12 <div id="nav">
13 <a href="#about">About</a>
14 <a href="#our-products">Our Products</a>
15 <a href="#contact">Contact</a>
16 <a href="#login">Login</a>
17 </div>
18 <div id="header">
19 <div id="content">
20 <div class="banner" style="background-image: url(images/banner1.jpg); background-size: 100%; background-size: cover">
21 <div class="text">Hello, my name is Chris!</div>
22 <div>This is the first banner at top of site</div>
23 </div>
24 <div class="container">
25 <div class="row">
26 <div class="col-4">
27 <div class="text">
28 <p>Nullam accumsan purus vel elit ultricies elementum. Vestibulum metus vehicula molestie. Nam porttitor magna id ligula tincidunt, id aliquam lacus tempus. Quam nunc tortor, ultrices ac est vitae, ultrices tempus magna. Pellentesque et rhoncus dui, sed lobortis diam. Ut eros ac nulla nec cursus. Proin laoreet eros arcu quis tincidunt.</p>
29 </div>
30 <div>This content will take up 3/12 (or 1/4) of the container</div>
31 </div>
32 <div class="col-4">
33 <div class="text">
34 
35 </div>
36 </div>
37 </div>
38 </div>
39 <div id="footer">
40 <div class="text">
41 <p>Footer copyright</p>
42 </div>
43 <div class="text">
44 <p>Contact footer</p>
45 </div>
46 </div>
```