

# Regressão Linear Múltipla

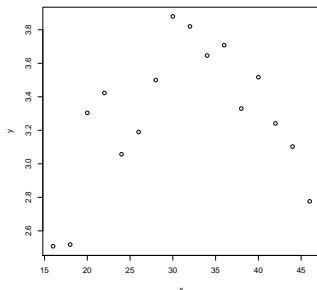
Prof. Dr. Leandro Balby Marinho



Aprendizagem de Máquina

## Regressão Polinomial

Em muitos casos o grafo de dispersão sugere uma relação não linear entre  $x$  e  $y$ .

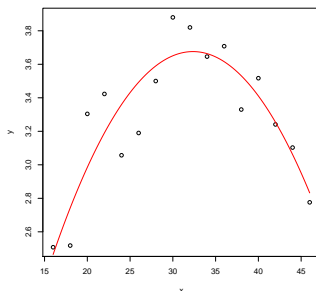


A equação do modelo quadrático, por exemplo, é dada por

$$y_i = w_0 + w_1x_i + w_2x_i^2 + \epsilon_i$$

# Regressão Polinomial

Em muitos casos o grafo de dispersão sugere uma relação não linear entre  $x$  e  $y$ .



A equação do modelo quadrático, por exemplo, é dada por

$$y_i = w_0 + w_1x_i + w_2x_i^2 + \epsilon_i$$

# Regressão Polinomial

Modelo:

$$y_i = w_0 + w_1x_i + w_2x_i^2 + \dots + w_px_i^p + \epsilon_i$$

- ▶ atributo 1 = 1 (constante)      parâmetro 1 =  $w_0$
- ▶ atributo 2 =  $x$       parâmetro 2 =  $w_1$
- ▶ atributo 3 =  $x^2$       parâmetro 3 =  $w_2$
- ▶ ...
- ▶ atributo  $p + 1 = x^p$       parâmetro  $p + 1 = w_d$

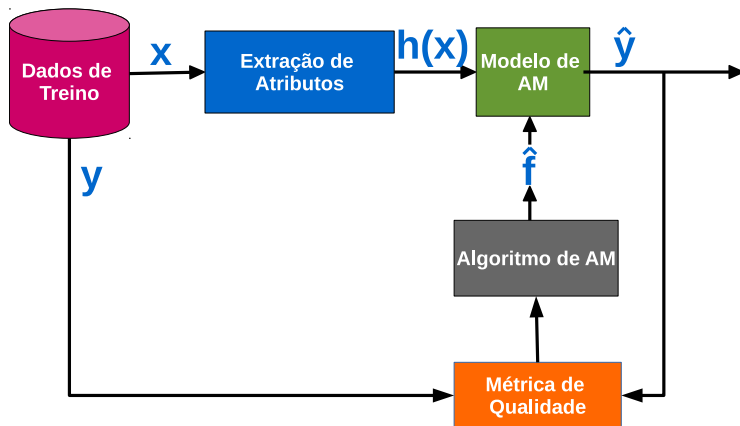
# Atributos como funções

Modelo:

$$\begin{aligned}y_i &= w_0 h_0(x_i) + w_1 h_1(x_i) + w_2 h_2(x_i) + \dots + w_p h_p(x_i) + \epsilon_i \\ &= \sum_{j=0}^D w_j h_j(x_i) + \epsilon_i\end{aligned}$$

- ▶ atributo 1 =  $h_0(x)$  ... geralmente 1 (constante)
- ▶ atributo 2 =  $h_1(x)$  ... e.g.,  $x$
- ▶ atributo 3 =  $h_2(x)$  ... e.g.,  $x^2$
- ▶ ...
- ▶ atributo  $p+1$  =  $h_p(x)$  ... e.g.,  $x^p$

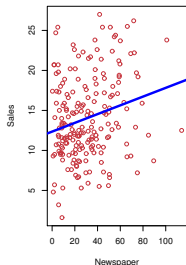
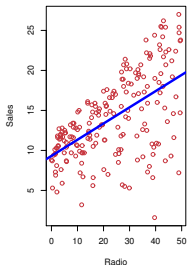
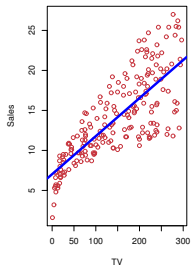
## Atributos como funções



# Roteiro

# Regressão Múltipla

Como usar as outras variáveis disponíveis no modelo de regressão?





# Notação

Saída:  $y$  (escalar)

Entradas:  $x = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d$

- ▶  $x^{(i)}$  =  $i$ -ésima observação (vetor)
- ▶  $x_j^{(i)}$  =  $j$ -ésima entrada da  $i$ -ésima observação (escalar)
- ▶  $h_j(x)$  =  $j$ -ésimo atributo (escalar)
- ▶ # observações  $(x, y)$ :  $N$
- ▶ # entradas  $x_j$ :  $d$
- ▶ # atributos  $h_j(x)$ :  $D$

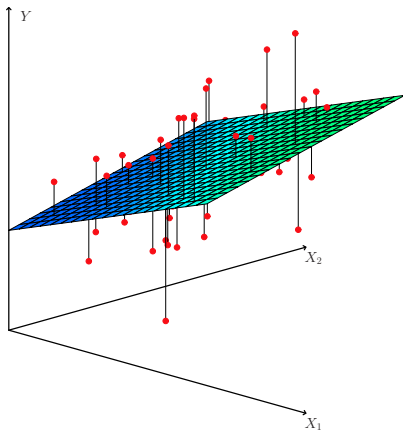
# Regressão Linear Múltipla

Modelo:

$$y^{(i)} = w_0 + w_1 x_1^{(i)} + w_2 x_2^{(i)} + \dots + w_d x_d^{(i)} + \epsilon_i$$

- ▶ variável 1 = 1
- ▶ variável 2 =  $x_1 \dots$  e.g., investimento em TV
- ▶ variável 3 =  $x_2 \dots$  e.g., investimento em Rádio
- ▶  $\dots$
- ▶ variável  $d + 1 = x_d \dots$  e.g., investimento em redes sociais

# Modelo de Regressão como Hiperplano



# Modelo de Regressão como uma Curva D-dimensional

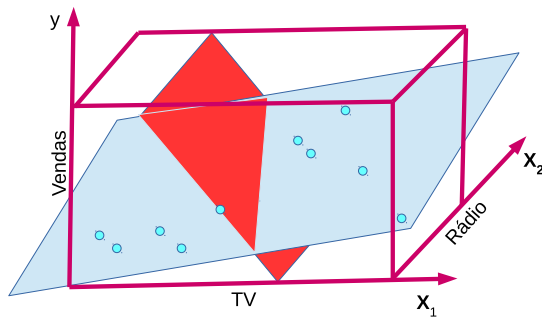
Modelo:

$$\begin{aligned}y^{(i)} &= w_0 + w_1 h_0(x^{(i)}) + w_2 h_1(x^{(i)}) + \dots + w_D h_D(x^{(i)}) + \epsilon_i \\ &= \sum_{j=0}^D w_j h_j(x^{(i)}) + \epsilon_i\end{aligned}$$

- ▶ atributo 1 =  $h_0(x)$  ... e.g., 1
- ▶ atributo 2 =  $h_1(x)$  ... e.g.,  $x_1$  = investimento em TV
- ▶ atributo 3 =  $h_2(x)$  ... e.g.,  $\log(x_2)x_1$
- ▶ ...
- ▶ atributo  $D + 1 = h_D(x)$  ... alguma outra função de  $x_1, \dots, x_D$

## Interpretando os coeficientes: dois atributos

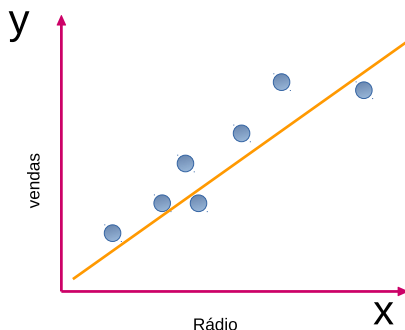
$$\hat{y} = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$



Fixando  $x_1$  a interpretação é a mesma da regressão linear simples.

## Interpretando os coeficientes: dois atributos

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x_1 + \hat{w}_2 x_2$$



Fixando  $x_1$  a interpretação é a mesma da regressão linear simples.

# Interpretando os coeficientes: múltiplos atributos

$$\hat{y} = \hat{w}_0 + \hat{w}_1 x_1 + \dots + \hat{w}_j x_j + \dots + \hat{w}_d x_d$$

Fixando todas as variáveis menos uma a interpretação é a mesma da regressão linear simples.

## Usando notação de vetores: uma observação

The diagram illustrates the vector notation for a single observation in multiple regression. It shows the equation  $y_i = \mathbf{w}^T \mathbf{h}(\mathbf{x}_i) + \varepsilon_i$ . On the left, a pink box contains  $y_i$ . This is followed by an equals sign. To the right of the equals sign is a horizontal row of six blue boxes, representing the feature vector  $\mathbf{h}(\mathbf{x}_i)$ . A blue bracket above this row is labeled  $\mathbf{w}^T$ , indicating the weight vector. To the right of the blue boxes is a vertical column of six green boxes, representing the error term  $\varepsilon_i$ . A blue bracket to the right of this column is labeled  $\mathbf{h}(\mathbf{x}_i)$ . A plus sign is placed between the green boxes and a grey box containing  $\varepsilon_i$ .

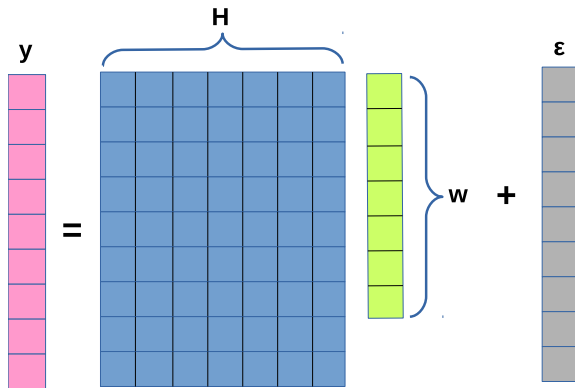


# Usando notação de vetores: uma observação

Para a observação  $i$ :

$$\begin{aligned}y^{(i)} &= \sum_{j=0}^D w_j h_j(x^{(i)}) + \epsilon^{(i)} \\ &= \mathbf{w}^T \mathbf{h}(x^{(i)}) + \epsilon^{(i)}\end{aligned}$$

Usando notação de matrizes: todas as observações



$$y = Hw + \epsilon$$

# Custo de uma curva D-dimensional

$$\begin{aligned}\text{RSS}(\mathbf{w}) &= \sum_{i=1}^N (y_i - h(\mathbf{x}^{(i)})^T \mathbf{w})^2 \\ &= (\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})\end{aligned}$$

# Gradiente do RSS

$$\begin{aligned}\nabla \text{RSS}(\mathbf{w}) &= \nabla [(\mathbf{y} - \mathbf{H}\mathbf{w})^T (\mathbf{y} - \mathbf{H}\mathbf{w})] \\ &= -2\mathbf{H}^T (\mathbf{y} - \mathbf{H}\mathbf{w})\end{aligned}$$

## Calculando parâmetros de forma fechada

$$\nabla \text{RSS}(\mathbf{w}) = -2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w}) = 0$$

Resolvendo para  $\mathbf{w}$ :

$$-2\mathbf{H}^T\mathbf{y} + 2\mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} = 0$$

$$-\mathbf{H}^T\mathbf{y} + \mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} = 0 \quad (\text{divide ambos os lados por 2})$$

$$\mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} = \mathbf{H}^T\mathbf{y}$$

$$(\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T\mathbf{H}\hat{\mathbf{w}} = (\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T\mathbf{y}$$

$$\hat{\mathbf{w}} = (\mathbf{H}^T\mathbf{H})^{-1}\mathbf{H}^T\mathbf{y}$$

Custo da inversão de matrizes (quando inversível):  $O(D^3)$

# Equações Normais como Sistemas de Equações Lineares

- ▶ As equações normais podem ser dadas por:

$$H^T H \hat{w} = H^T y$$

- ▶ A equação acima pode ser representada por um sistema de equações lineares da forma:

$$\underbrace{H^T H}_A \underbrace{w}_x = \underbrace{H^T y}_b$$

Vários métodos de resolução:

- ▶ Eliminação Gaussiana
- ▶ Fatoração de Cholesky
- ▶ Fatoração QR

# Exemplo 1

Use um modelo de regressão linear múltipla para estimar o valor de  $y$  para  $x_1 = 3$  e  $x_2 = 4$  considerando os dados abaixo.

$x_1$	$x_2$	$y$
1	2	3
2	3	2
4	1	7
5	5	1

## Exemplo 1 cont.

Modelo regressão múltipla:

$$y_i = w_0 + w_1x_1 + w_2x_2 + \epsilon_i$$

$$H = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 3 \\ 1 & 4 & 1 \\ 1 & 5 & 5 \end{bmatrix}, \quad y = \begin{bmatrix} 3 \\ 2 \\ 7 \\ 1 \end{bmatrix}$$

$$H^T H = \begin{bmatrix} 4 & 12 & 11 \\ 12 & 46 & 37 \\ 11 & 37 & 39 \end{bmatrix}, \quad H^T y = \begin{bmatrix} 13 \\ 40 \\ 24 \end{bmatrix}$$



## Exemplo 1 cont.

Estimando os parâmetros por Eliminação Gaussiana:

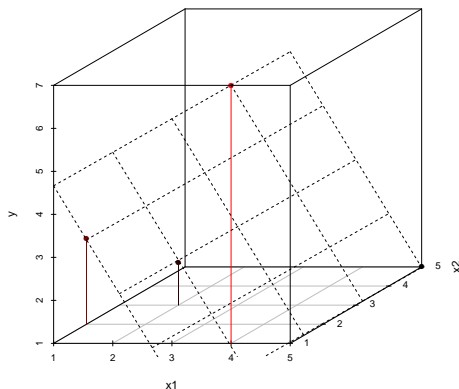
$$\left[ \begin{array}{ccc|c} 4 & 12 & 11 & 13 \\ 12 & 46 & 37 & 40 \\ 11 & 37 & 39 & 24 \end{array} \right] \rightarrow \left[ \begin{array}{ccc|c} 1 & 3 & 2.75 & 3.25 \\ 0 & 10 & 4 & 1 \\ 0 & 4 & 8.75 & -11.75 \end{array} \right]$$

$$\rightarrow \left[ \begin{array}{ccc|c} 4 & 12 & 11 & 13 \\ 0 & 1 & 0.4 & 0.1 \\ 0 & 0 & 7.15 & -12.15 \end{array} \right]$$

$$w \approx \left[ \begin{array}{c} 5.583 \\ 0.779 \\ -1.699 \end{array} \right]$$

## Exemplo 1 cont.

$$\hat{y}(x_1 = 3, x_2 = 4) = 5.583 + 0.779x_1 - 1.699x_2 = 1.124$$



# Gradiente Descendente

## Gradient-Descent

1 **while** not converged

2 
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \alpha \underbrace{\nabla \text{RSS}(\mathbf{w}^{(t)})}_{-2\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{w})}$$

# Gradiente Descendente

## Gradient-Descent

1 **while** not converged

2 
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + 2\alpha \mathbf{H}^T (y - \underbrace{\mathbf{H}\mathbf{w}^{(t)}}_{\hat{y}})$$

## Derivada parcial de um parâmetro

$$\begin{aligned}
 \text{RSS}(\mathbf{w}) &= \sum_{i=1}^N (y_i - \mathbf{h}(\mathbf{x}^{(i)})^T \mathbf{w})^2 \\
 &= \sum_{i=1}^N (y^{(i)} - w_0 h_0(\mathbf{x}^{(i)}) - w_1 h_1(\mathbf{x}^{(i)}) - \dots - w_D h_D(\mathbf{x}^{(i)}))
 \end{aligned}$$

Derivada parcial em relação a  $w_j$

$$\begin{aligned}
 &= \sum_{i=1}^N 2(y^{(i)} - w_0 h_0(\mathbf{x}^{(i)}) - w_1 h_1(\mathbf{x}^{(i)}) - \dots - w_D h_D(\mathbf{x}^{(i)}))(-h_j(\mathbf{x}^{(i)})) \\
 &= -2 \sum_{i=1}^N h_j(\mathbf{x}^{(i)})(y^{(i)} - \mathbf{h}(\mathbf{x}^{(i)})^T \mathbf{w})
 \end{aligned}$$

# Algoritmo do Gradiente Descendente

GradientDescent( $\alpha, \epsilon$ )

```
1  initialize  $\mathbf{w}$ ,  $t = 1$ 
2  while  $\|\nabla \text{RSS}(\mathbf{w}^{(t)})\| \geq \epsilon$ 
3      for  $j = 0, \dots, D$ 
4           $\text{partial}[j] = -2 \sum_{i=1}^N h_j(\mathbf{x}^{(i)})(y^{(i)} - h(\mathbf{x}^{(i)})^T \mathbf{w})$ 
5           $\mathbf{w}_j^{(t+1)} = \mathbf{w}_j^{(t)} - \alpha \cdot \text{partial}[j]$ 
6       $t = t + 1$ 
7  return  $\mathbf{w}$ 
```

# Gradiente Descendente vs. Equações Normais

- ▶ Gradiente Descendente
  - ▶ Precisa escolher  $\alpha$ .
  - ▶ Pode precisar de muitas iterações.
  - ▶ Relativamente eficiente para  $D$  grande.
- ▶ Equações Normais
  - ▶ Não precisa escolher  $\alpha$ .
  - ▶ Não precisa iterar.
  - ▶ Métodos de resolução de sistemas de equações lineares podem ser caros (e.g. fatoração de Cholesky  $\in O(D^3)$ ).
  - ▶ Lento para  $D$  muito grande.

# Referências



Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani. An Introduction to Statistical Learning with Applications in R. Springer, 2013.



Emily Fox and Carlos Guestrin. Machine Learning Specialization. Curso online disponível em <https://www.coursera.org/specializations/machine-learning>  
Último acesso: 04/09/2017.