# Introdução ao Aprendizado de Máquina para Físicos

Marcelo Vargas dos Santos
Aula 2

# O ovo ou a galinha?



Qual caminho seguir?

1. Formular uma pergunta a ser respondida com os dados que temos?

2. Ou buscar dados para responder uma pergunta já formulada?

Tanto faz! Ambos os caminhos são válidos.
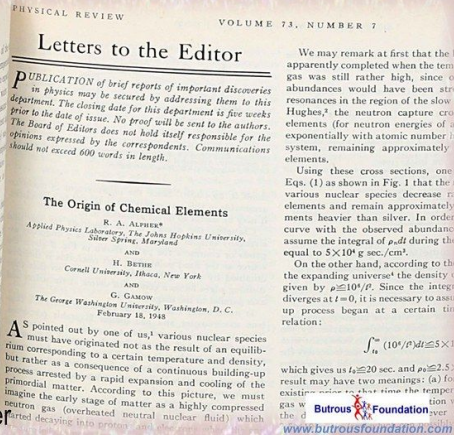
# Radiação Cósmica de Fundo



**The Alpher–Bethe–Gamow paper,** or αβγ **paper,** was published in Physical Review on 1st April 1948 by the graduate student Ralph Alpher, and his advisor George Gamow. *The work, argued that the Big Bang would create hydrogen, helium and heavier elements in the correct proportions to explain their abundance in the early universe.*

*Their work affirmed that the extreme conditions at the start-up of the universe could explain the existing abundance of its most common elements.*

George Gamow added Bethe's name (in absentia) without consulting him, knowing that Bethe would not mind, and against Ralph Alpher's wishes. This was apparently a reflection of Gamow's sense of humor, wanting to have a paper title that would sound like *the first three letters* of the Greek alphabet. As one of the Physical Review's reviewers, Bethe saw the manuscript and struck out the words "in absentia"

Ralph Alpher α  Hans Albrecht Bethe β  George Gamow γ  paper

www.butrousfoundation.com

# Onde encontrar dados

1. Google Dataset Dearch:datasetsearch.research.google.com
2. Kaggle: kaggle.com (Competições em ciência de dados)
3. Drivendata: drivendata.org (Competições)
4. Portal Brasileiro de Dados: dados.gov.br
5. 538: fivethirtyeight.com (Opinião Pública)
6. Quandl: quandl.com (Dados Financeiros)
7. Reddit: reddit.com/r/datasets

# E depois? Análise Exploratória de dados





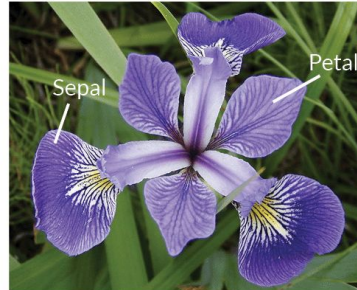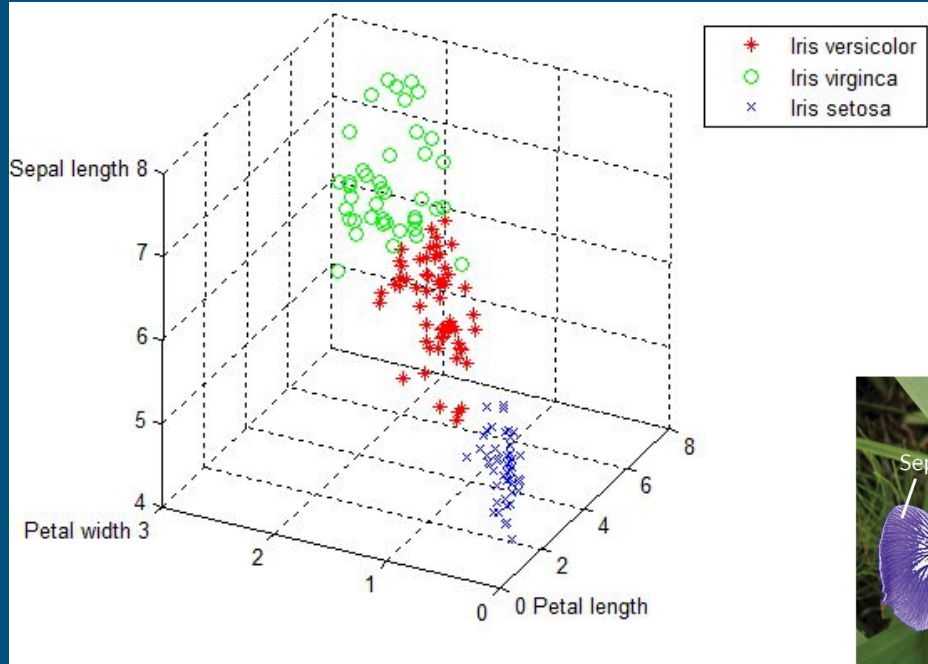A comprehensive review of tools for exploratory analysis of tabular industrial datasets
https://doi.org/10.1016/j.visinf.2018.12.004

Escolhendo o método

# Não supervisonado (Agrupamento)



Iris Versicolor
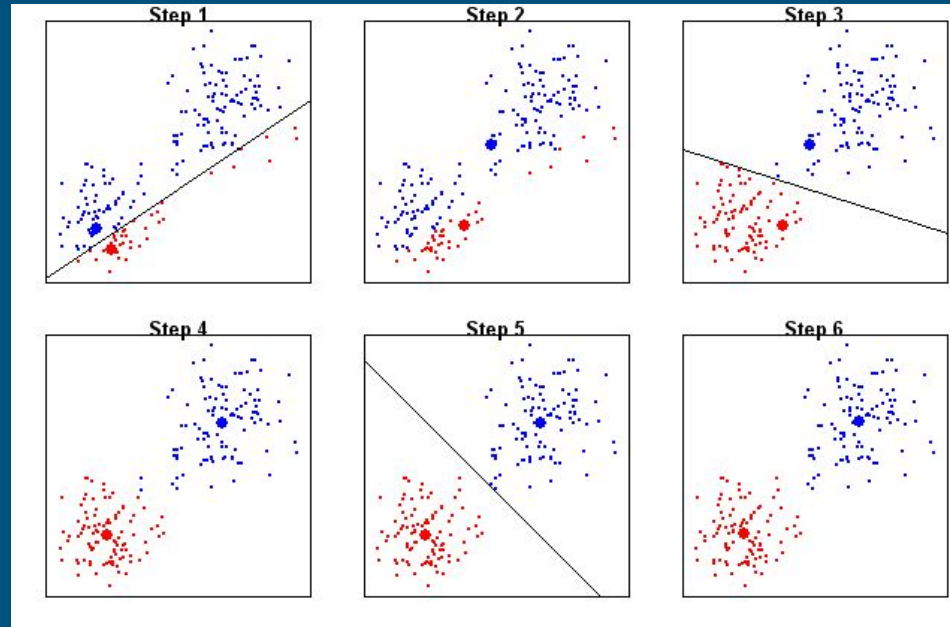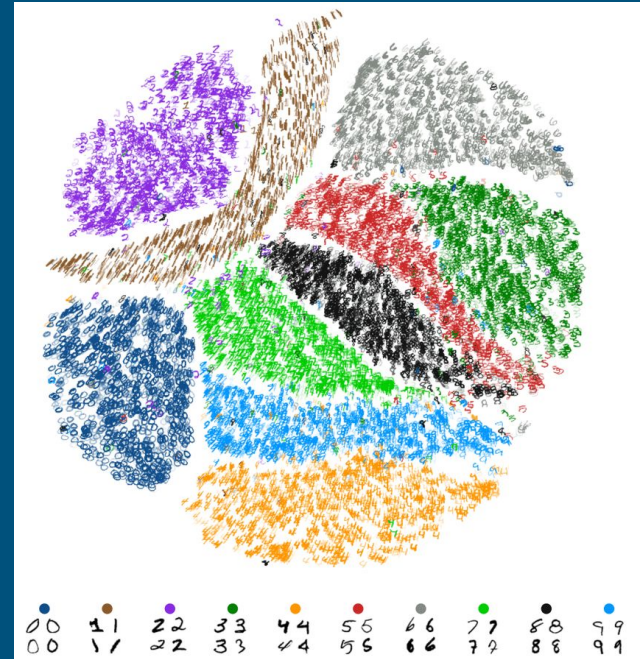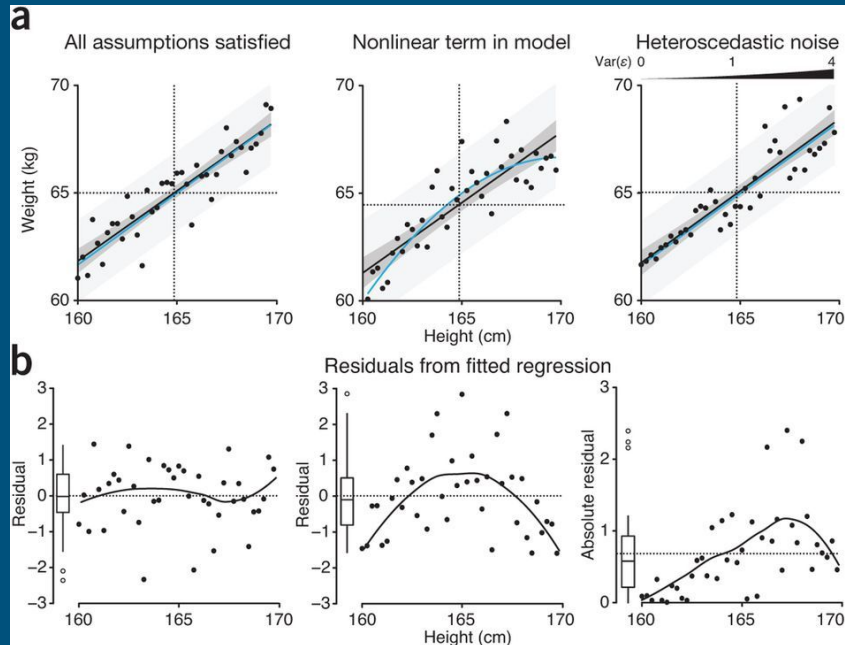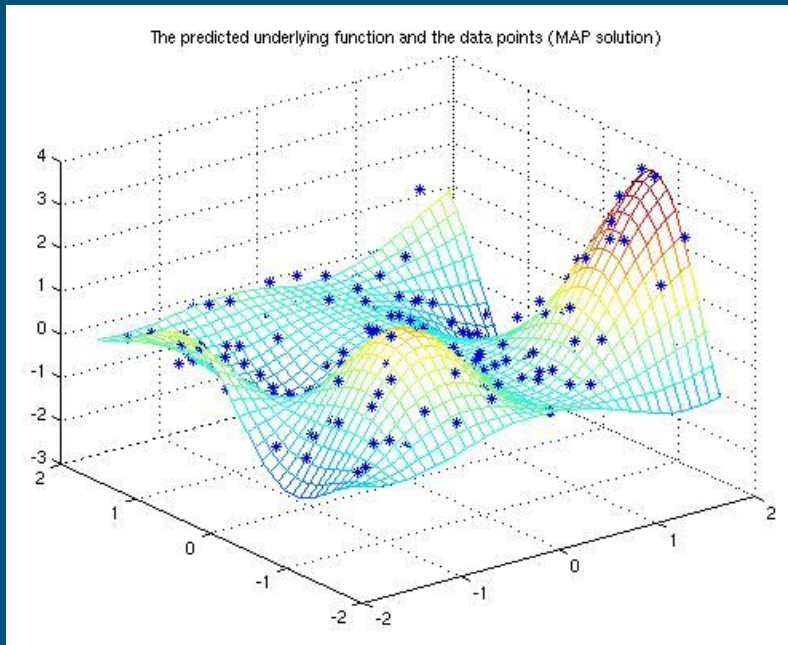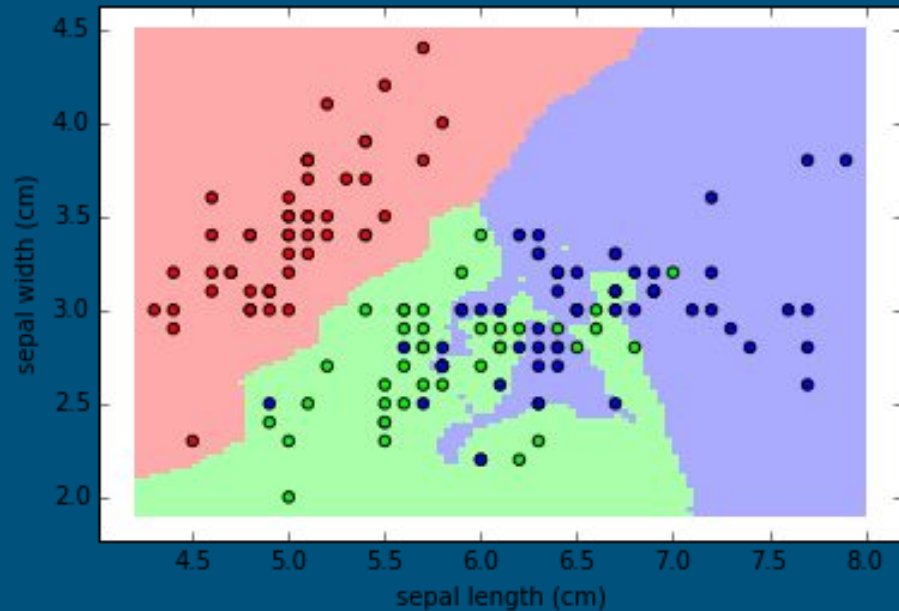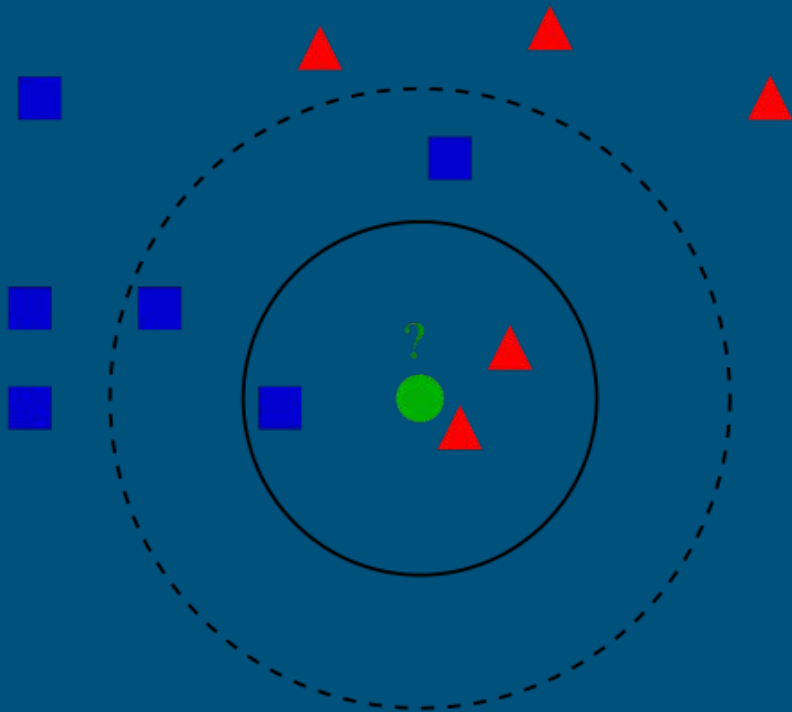
Iris Setosa

Iris Virginica

# Modelo K-Means

# Classificação

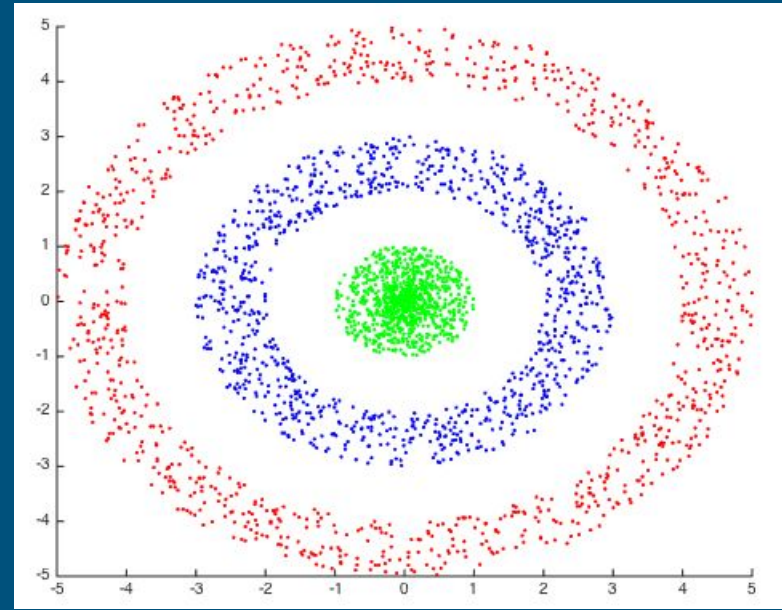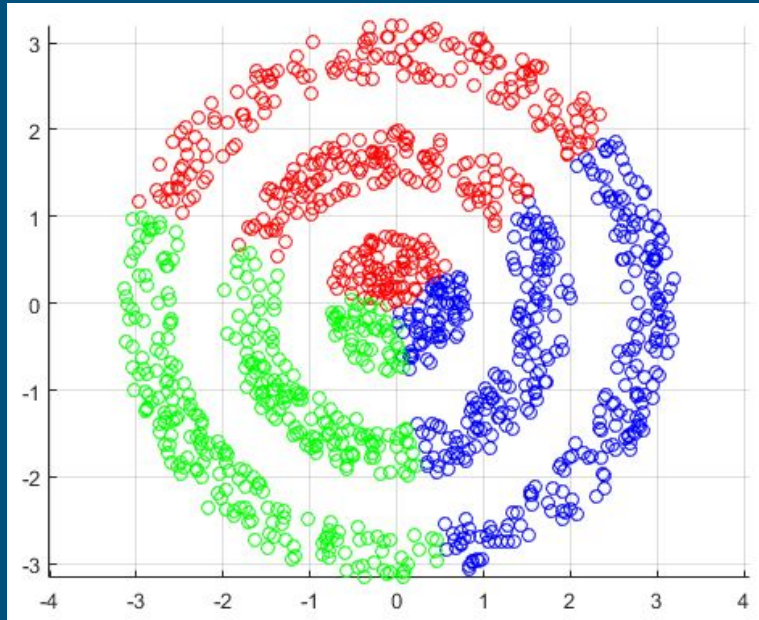# Regressão

# K-ésimo Vizinho mais Próximo
## k-nearest neighbors

# Pré Processamento

# Processo de Aprendizado

# Ingredientes

$$\mathbf{X} \qquad : \qquad \text{Amostra}$$

$$\mathcal{G}(\cdot,\cdot) \qquad : \qquad \text{Algoritmo}$$

$$g_{\mathbf{w},\mathbf{X}}(\cdot) = \mathcal{G}(\mathbf{w},\mathbf{X}) \qquad : \qquad \text{Modelo}$$

$$\mathbf{w} \qquad : \qquad \text{Hiperparâmetros}$$

$$\mathbf{Y} = g_{\mathbf{w},\mathbf{X}}(\bar{\mathbf{X}}) \qquad : \qquad \text{Previsão}$$

$$\mathcal{C}(\bar{\mathbf{X}},\mathbf{Y}) \qquad : \qquad \text{Custo}$$

# Processo iterativo

## Pseudo-Código

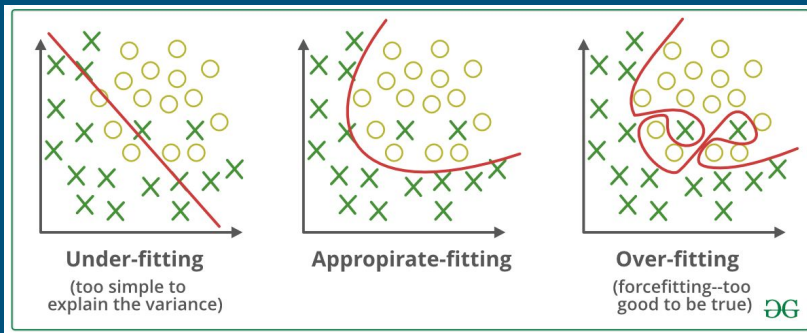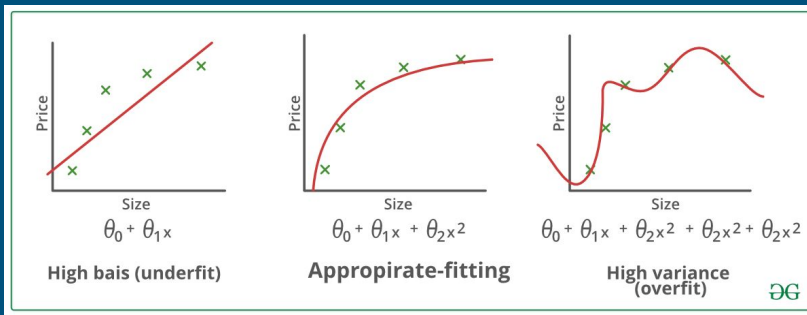$$\overline{\mathbf{X}} = \mathbf{X}_{Test} + \mathbf{X}_{Train}$$

Para $\mathbf{W}$ em $(\mathbf{w}_1, \mathbf{w}_2, \ldots)$ :

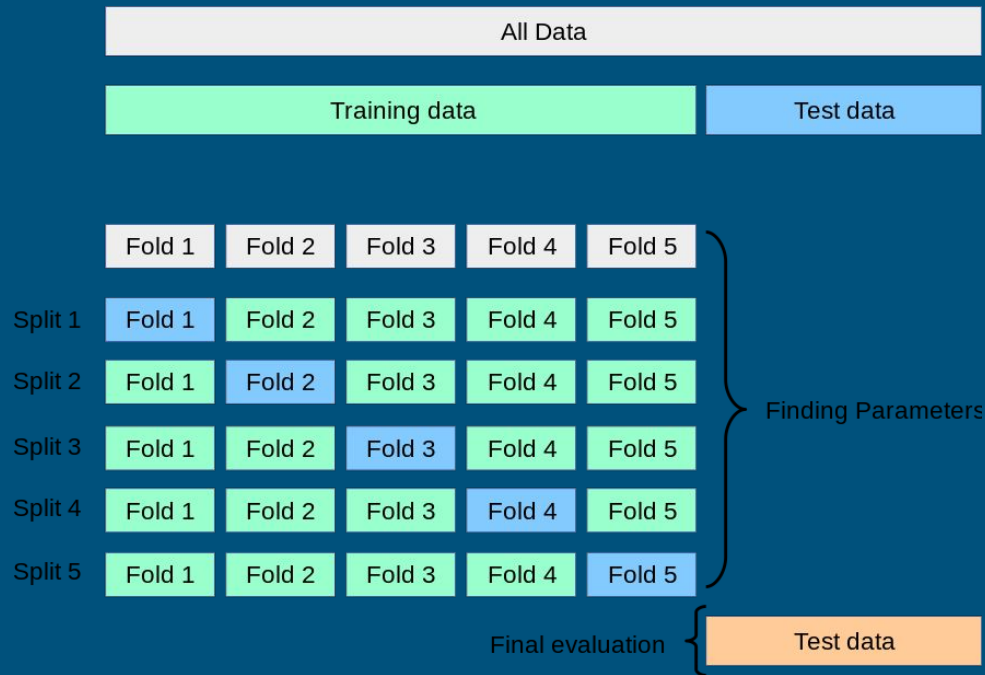$$g = \mathcal{G}(\mathbf{w}, \mathbf{X}_{Train})$$

$$c_i = \mathcal{C}(\mathbf{X}_{Test}, g(\mathbf{X}_{Test})$$

# Super-ajuste vs Sub-ajuste Overfitting vs Underfitting

# Validação cruzada

## Pseudo-Código

$$\mathbf{X} = \mathbf{X}_{Test} + \mathbf{X}_{Train}$$

Para $\mathbf{W}$ em $(\mathbf{w}_1, \mathbf{w}_2, \ldots)$ :

$$g = \mathrm{CV}(\mathcal{G}(\mathbf{w}, \cdot), \mathbf{X}_{Train})$$

$$c_i = \mathcal{C}(\mathbf{X}_{Test}, g(\mathbf{X}_{Test})$$

# Análise Exploratória de Dados

Próxima Aula