



KNIME

Machine learning para não programadores

GLEBER TACIO TEIXEIRA, DSC.

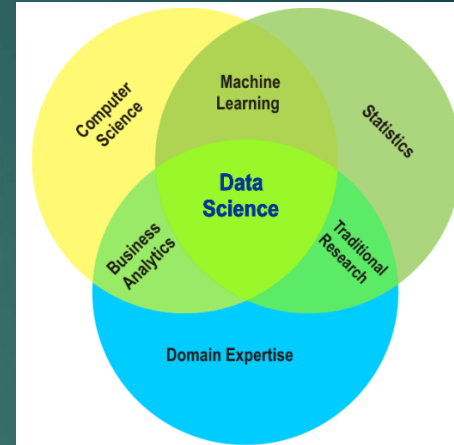
PETROBRAS

Agenda

- ▶ Introdução
- ▶ Algoritmos de Machine Learning
 - ▶ Linear regression
 - ▶ Decision tree
 - ▶ Random forest
 - ▶ MLP
- ▶ Ferramenta Knime
- ▶ Conclusões

Introdução

- ▶ O que é Machine Learning?



<http://www.marsiantech.com/best-data-science-training-in-pune.php>

Machine Learning = Computação + Estatística + Álgebra



Modelos

Introdução

- ▶ É possível trabalhar com *Machine Learning* sem ser matemático, estatístico ou da área de computação?



Geoffrey E. Hinton

<http://www.cs.toronto.edu/~hinton/geoff6.jpg>

Introdução

- ▶ É possível trabalhar com *Machine Learning* sem gostar de matemática?



Aprender
matemática você
deve, jovem
Padawan!

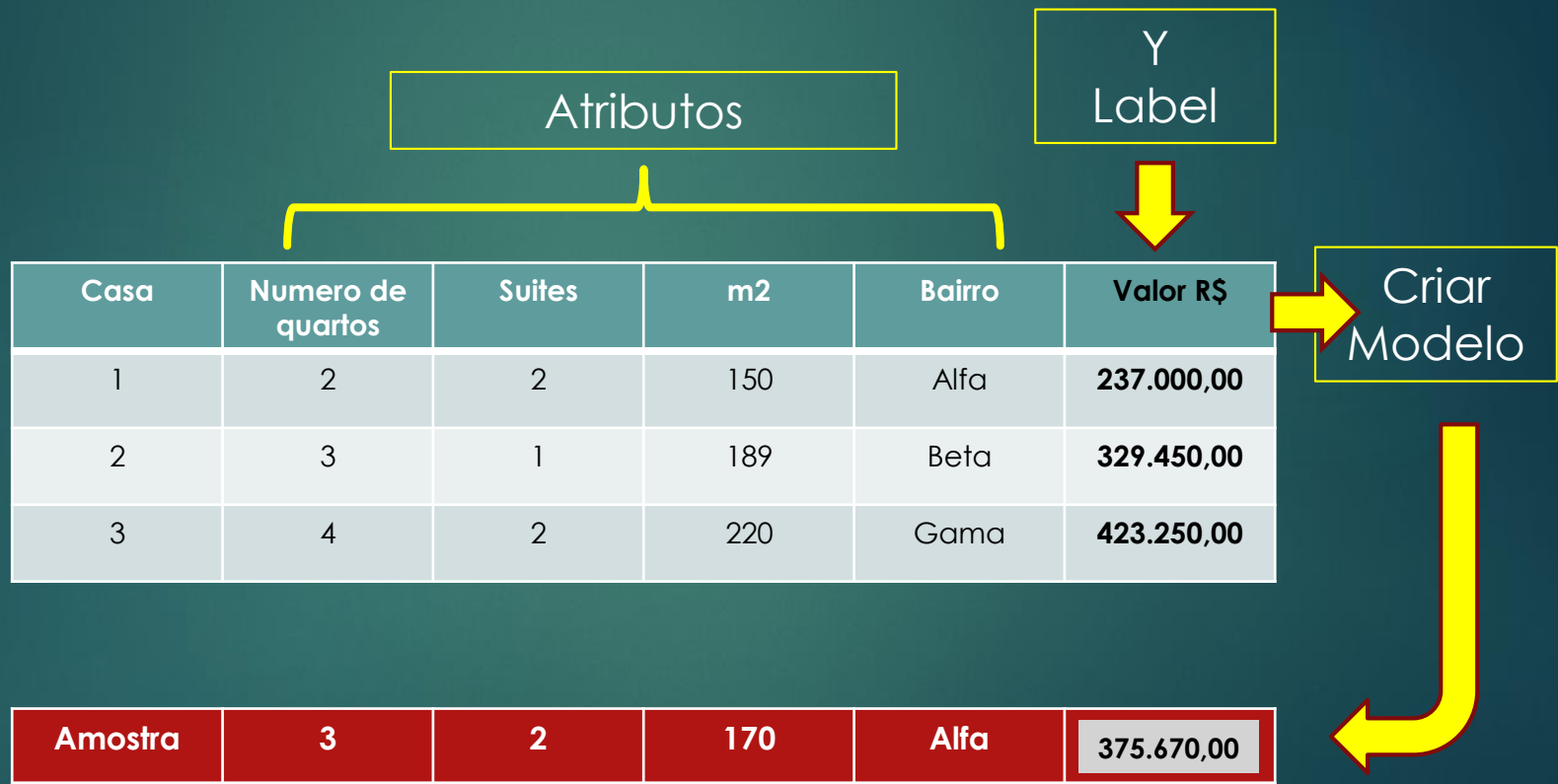
Introdução

- ▶ É possível trabalhar com *Machine Learning* sem saber programar?



Algoritmos de *Machine Learning*

Método Supervisionado - regressão



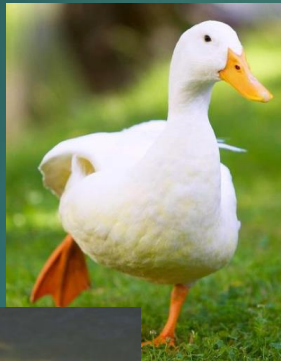
Algoritmos de Machine Learning

Supervisionado - classificação

Pato



Pato



???

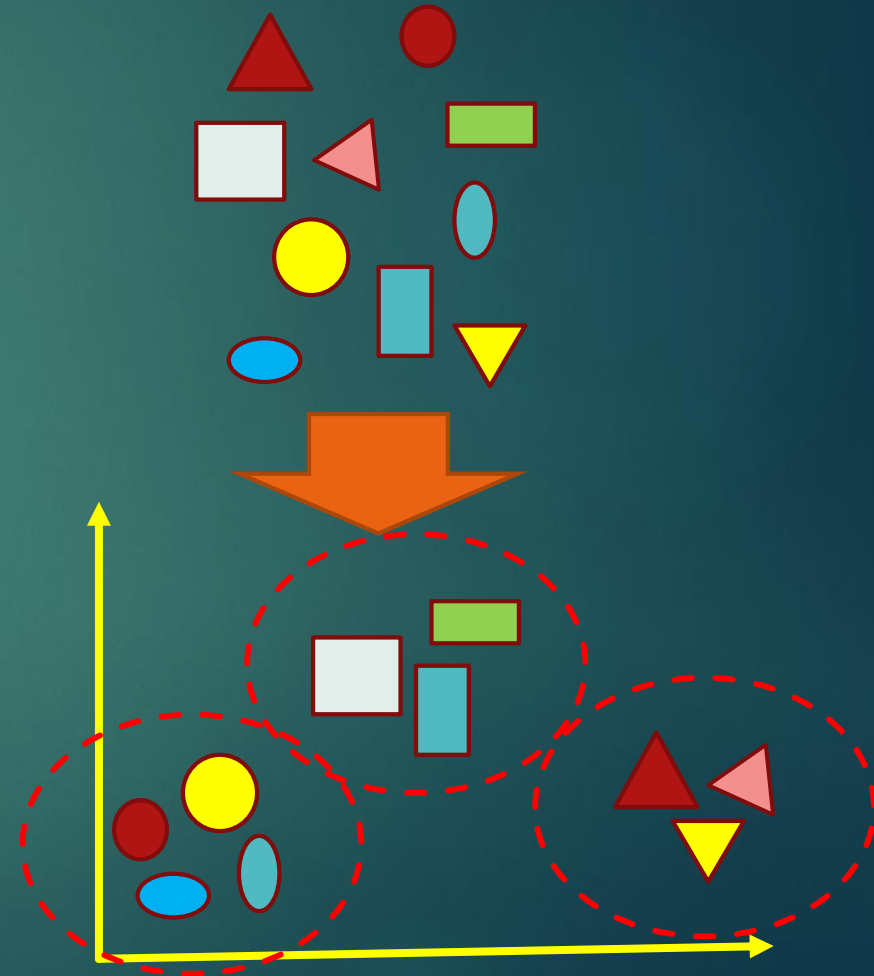


Castor



Castor

Não Supervisionado



Algoritmos de *Machine Learning*

Variáveis contínuas

Não Supervisionado

Clusterização e redução
de dimensionalidade
K-means
PCA

Supervisionado (Regressão)

Regressão Linear
Decision Tree
GB Tree
Random Forests
MLP

Variáveis categóricas

Não Supervisionado

Análise de associação
Markov Models

Supervisionado (Classificação)

Decision Tree
Random Forests
MLP
Regressão Logística
Naive-Bayes
SVM
KNN

Terminologia na modelagem preditiva

- ▶ Modelo: função explicativa.
- ▶ Dados: o que você está usando pra encaixar no modelo.
- ▶ Alvo/Target/Label/Variável dependente/ y : valor que você está tentando prever.
- ▶ Características/Features/Variáveis independentes/ X atributos que serão utilizados na predição.
- ▶ Métodos: algoritmos que usarão seus dados para obter um modelo.

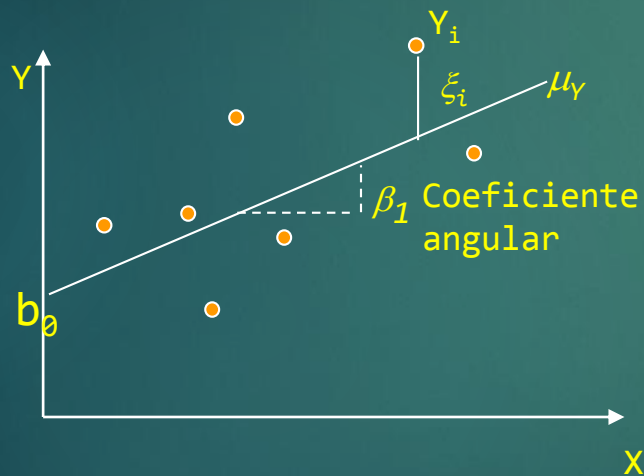
Regressão linear

Diagram illustrating the components of the linear regression equation:

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

Labels and arrows pointing to the equation components:

- Variável Dependente** points to Y_i .
- Intercepto** points to β_0 .
- Inclinação** points to β_1 .
- Variável Independente** points to X_i .
- Erro Aleatório** points to ε_i .



$$\hat{Y}_i = b_0 + b_1 X_i \quad \text{Modelo estimado}$$

$$\varepsilon_i = Y_i - \hat{Y}_i \quad \text{Resíduo}$$

Vantagens

- Fácil de avaliar a qualidade do modelo
- Bem fundamentado
- Intuitivo

Desvantagens

- Cuidado ao se trabalhar com variáveis categóricas (qualitativa)
- Deve-se cumprir os pressupostos (Linearidade, homocedasticidade....)

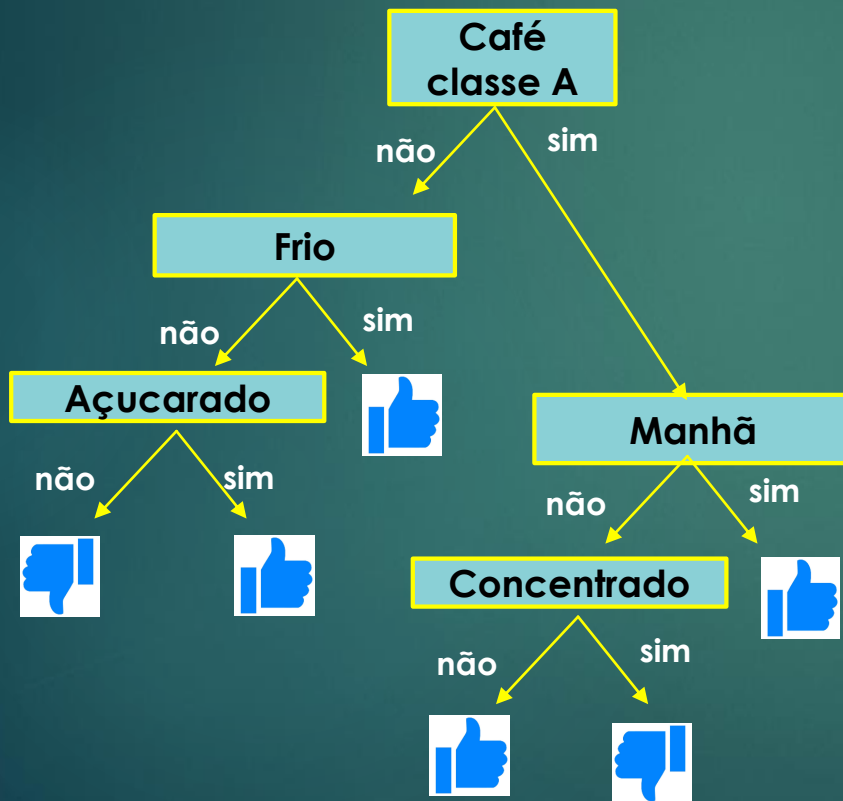


Classe A	Dia Frio	Manhã	Concentrado	Açucarado	Resultado
Sim	Sim	Não	Não	Sim	
Não	Sim	Sim	Sim	Não	
Não	Não	Sim	Sim	Sim	
Sim	Nao	Não	sim	Sim	

...

Decision Tree

- ▶ “Aprende” criando ramos ou subdivisões do dataset de maneira recursiva com ganho de informação (supervisionado)



Vantagens

- Simple de usar e interpretar (gráfico da árvore)
- Robustas
- Variáveis categorias ou numéricas
- Pouca preparação
- Trabalha bem com muitos dados

Desvantagens

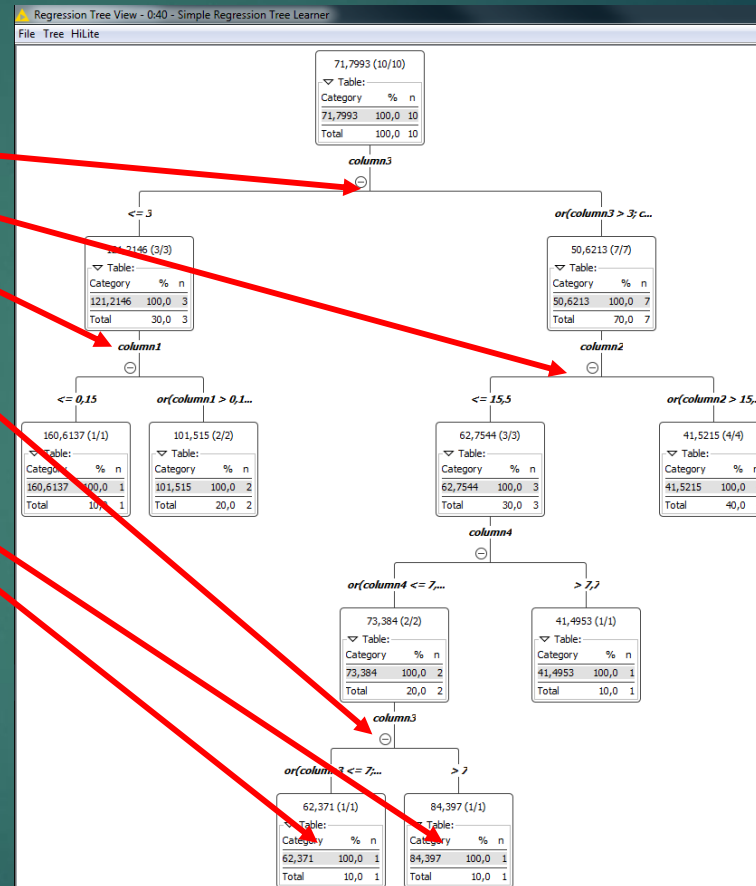
- Podem não ser muito precisas
- Podem sofrer overfitting.
- Variáveis categóricas com várias divisões podem induzir desvio.

Árvore de decisão

Número de nós

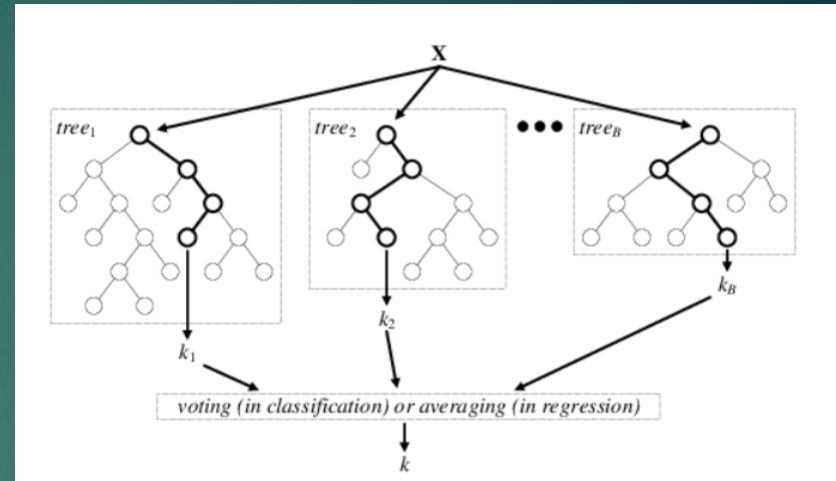
Número de divisões de
nós
(filhos)

Prof. da árvore, 4
camadas



Random forest

- ▶ Bom pra regressão e classificação
- ▶ Visualização do modelo
- ▶ Menor risco de overfitting
- ▶ Exige pouco preparo das amostras
- ▶ Exige pouco ajuste de hyperparâmetros
- ▶ Rápido pra treinar mas lento pra usar (real time)
- ▶ Bagging de árvores (decision tree)
 - ▶ Combinação de vários modelos pode gerar um modelo melhor
 - ▶ Votação

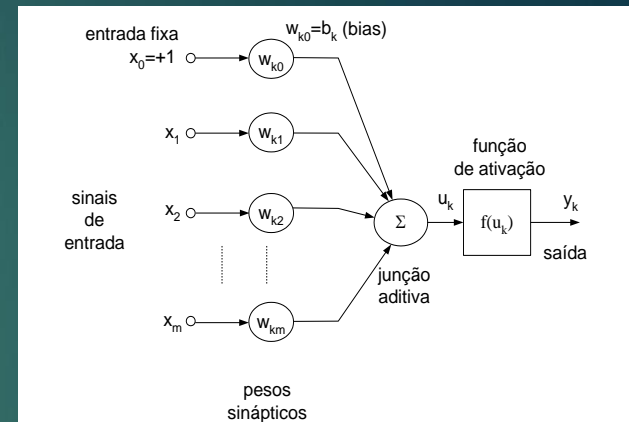


https://www.researchgate.net/figure/Architecture-of-the-random-forest-model_fig1_301638643

Rede Neural (Rprop - MLP)

Técnica inspirada no funcionamento do cérebro, onde neurônios artificiais, conectados em rede, são capazes de aprender e de generalizar.

Técnica de aproximação de funções por regressão não linear.



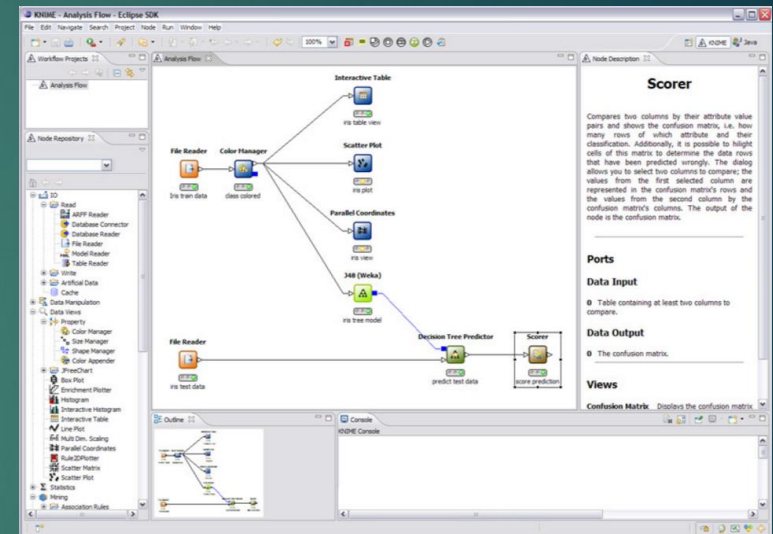


KNIME Analytics Platform

KNIME® Analytics Platform is the leading open solution for data-driven innovation, helping you discover the potential hidden in your data, mine for fresh insights, or predict new futures. Our enterprise-grade, open source platform is fast to deploy, easy to scale and intuitive to learn.

With more than 1000 modules, hundreds of ready-to-run examples, a comprehensive range of integrated tools, and the widest choice of advanced algorithms available, KNIME Analytics Platform is the perfect toolbox for any data scientist. Our steady course on unrestricted open source is your passport to a global community of data scientists, their expertise, and their active contributions.

- Free
- Intuitiva
- Fácil
- Ágil
- Open Source
- Constante Atualização
- Ampla aplicação



www.knime.com

Knime - Konstanz Information Miner

- ▶ Universidade de Konstanz – início em 2004 (Michael Berthold)
 - ▶ Proposta de plataforma de integração
- ▶ Primeira versão 2006
 - ▶ Companhias farmacêuticas e de softwares científicos
- ▶ Hoje - Plataforma de integração de vários projetos / empresas
- ▶ Desenvolvimentos atuais:
 - ▶ Big Data
 - ▶ Spark

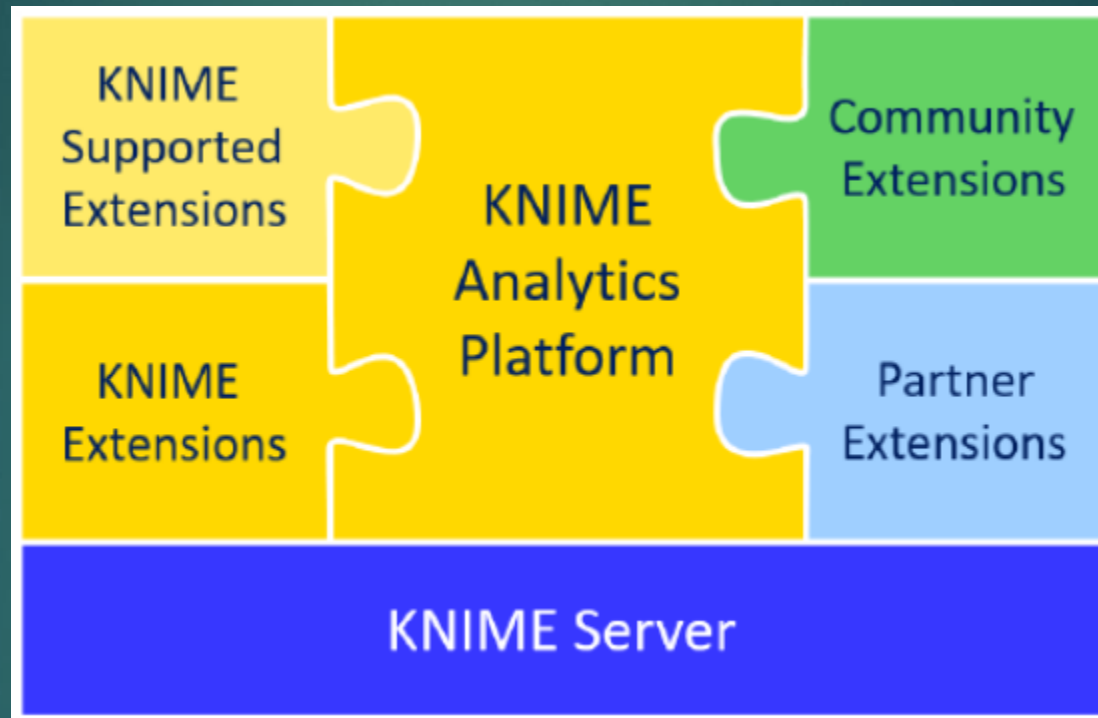
Knime - Konstanz Information Miner

- ▶ Plataforma analítica
 - ▶ Machine Learning
 - ▶ Data mining
 - ▶ ETL
- ▶ Plataforma de integração e de reporting.
- ▶ Mínima programação
 - ▶ Aplicação final
 - ▶ Prototipagem
 - ▶ Ensino



Fig. 1: Gartner 2018 Magic Quadrant for Data Science and Machine Learning Platforms

Knime - Konstanz Information Miner

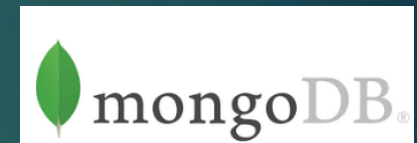


www.knime.com

Knime - Konstanz Information Miner



Perl



Ao trabalho ... Regressão

https://archive.ics.uci.edu/ml/datasets.html?format=&task=reg&att=&area=&numAtt=&numIns=&type=&sort=nameUp&view=table tradutor google 2017

UCI Machine Learning Repository
Center for Machine Learning and Intelligent Systems








About Citation Policy Donate a Data Set Contact

Search

Repository Web Google

[View ALL Data Sets](#)

Browse Through: **86 Data Sets** [Table View](#) [List View](#)

Default Task - Undo	Name	Data Types	Default Task	Attribute Types	# Instances	# Attributes	Year
Classification (322) Regression (86) Clustering (77) Other (54)	 3D Road Network (North Jutland, Denmark)	Sequential, Text	Regression, Clustering	Real	434874	4	2013
Attribute Type	 Air Quality	Multivariate, Time-Series	Regression	Real	9358	15	2016
Categorical (1) Numerical (78) Mixed (5)	 Air quality	Multivariate, Time-Series	Regression	Real	9358	15	2016
Data Type	 Airfoil Self-Noise	Multivariate	Regression	Real	1503	6	2014
Multivariate (76) Univariate (6) Sequential (7) Time-Series (26) Text (7) Domain-Theory (4) Other (0)	 Amazon Access Samples	Time-Series, Domain-Theory	Regression, Clustering, Causal-Discovery		30000	20000	2011
Area	 Appliances energy prediction	Multivariate, Time-Series	Regression	Real	19735	29	2017
Life Sciences (10) Physical Sciences (10) CS / Engineering (41) Social Sciences (7) Business (9) Game (1) Other (8)	 Auto MPG	Multivariate	Regression	Categorical, Real	398	8	1993
# Attributes							

Regressão – Resistência compressiva do concreto

UCI



Machine Learning Repository

Center for Machine Learning and Intelligent Systems

[About](#) [Citation Policy](#)

Repository

Concrete Compressive Strength Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Concrete is the most important material in civil engineering. The concrete compressive strength is a highly nonlinear function of age and ingredients.



Data Set Characteristics:	Multivariate	Number of Instances:	1030	Area:	Physical
Attribute Characteristics:	Real	Number of Attributes:	9	Date Donated	2007-08-03
Associated Tasks:	Regression	Missing Values?	N/A	Number of Web Hits:	130670

Source:

Original Owner and Donor
Prof. I-Cheng Yeh
Department of Information Management
Chung-Hua University,
Hsin Chu, Taiwan 30067, R.O.C.
e-mail: icyeh '@' chu.edu.tw
TEL: 886-3-5186511

Date Donated: August 3, 2007

Attribute Information:

Given are the variable name, variable type, the measurement unit and a brief description. The concrete compressive strength is the output variable of the database.

Name -- Data Type -- Measurement -- Description

Cement (component 1) -- quantitative -- kg in a m3 mixture -- Input Variable
Blast Furnace Slag (component 2) -- quantitative -- kg in a m3 mixture -- Input Variable
Fly Ash (component 3) -- quantitative -- kg in a m3 mixture -- Input Variable
Water (component 4) -- quantitative -- kg in a m3 mixture -- Input Variable
Superplasticizer (component 5) -- quantitative -- kg in a m3 mixture -- Input Variable
Coarse Aggregate (component 6) -- quantitative -- kg in a m3 mixture -- Input Variable
Fine Aggregate (component 7) -- quantitative -- kg in a m3 mixture -- Input Variable
Age -- quantitative -- Day (1~365) -- Input Variable
Concrete compressive strength -- quantitative -- MPa -- Output Variable

Classificação – Iris flower dataset



Ronald Fisher

Attribute Information:

1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm
5. class:
 - Iris Setosa
 - Iris Versicolour
 - Iris Virginica

UCI



Machine Learning Repository

Center for Machine Learning and Intelligent Systems

Iris Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Famous database; from Fisher, 1936



Data Set Characteristics:	Multivariate	Number of Instances:	150	Area:	Life
Attribute Characteristics:	Real	Number of Attributes:	4	Date Donated	1988-07-01
Associated Tasks:	Classification	Missing Values?	No	Number of Web Hits:	1954925

Source:

Creator:

R.A. Fisher

Donor:

Michael Marshall (MARSHALL%PLU"@io.arc.nasa.gov)

Download: [Data Folder](#), [Data Set Description](#)

Source:

Attribute Information:

- 1 - fixed acidity
- 2 - volatile acidity
- 3 - citric acid
- 4 - residual sugar
- 5 - chlorides
- 6 - free sulfur dioxide
- 7 - total sulfur dioxide
- 8 - density
- 9 - pH
- 10 - sulphates
- 11 - alcohol

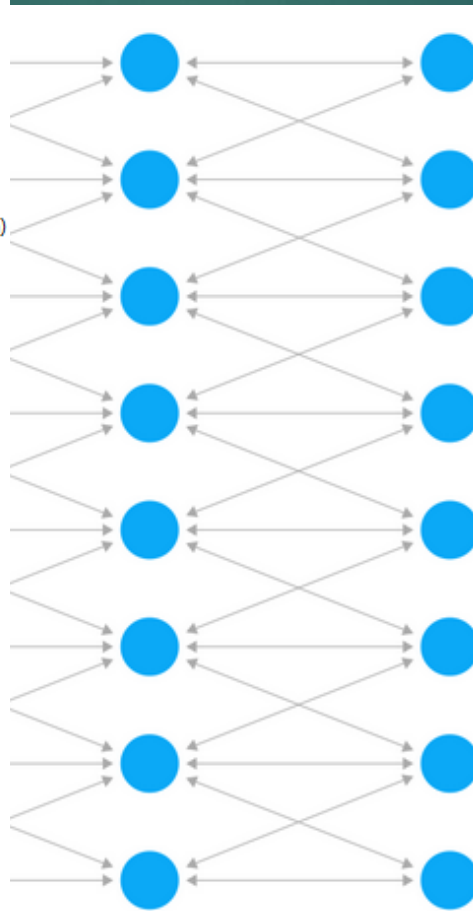
Output variable (base wine quality):

- 12 - quality (score between 4 and 10)

The two datasets are related to red and white variants of the Portuguese "Vinho Verde" wine. For more details, consult: [\[Web Link\]](#) or the reference [Cortez et al., 2009]. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g. there is no data about grape types, wine brand, wine selling price, etc.).

These datasets can be viewed as classification or regression tasks. The classes are ordered and not balanced (e.g. there are much more normal wines than excellent or poor ones). Outlier detection algorithms could be used to detect the few excellent or poor wines. Also, we are not sure if all input variables are relevant. So it could be interesting to test feature selection methods.

- KNIME Labs
 - Deep Learning
 - DL4J
 - I/O
 - DL4J Model Reader
 - DL4J Model Writer
 - Layer
 - Convolutional
 - Convolution Layer
 - LRN Layer
 - Pooling Layer
 - Encoder
 - AutoEncoder Layer
 - Perceptron
 - Dense Layer
 - DL4J Model Initializer
 - Learn
 - Supervised
 - DL4J Feedforward Learner (Classification)
 - DL4J Feedforward Learner (Regression)
 - Unsupervised
 - DL4J Feedforward Learner (Pretraining)
 - Networks
 - AlexNet
 - DeepBelief
 - DeepMLP
 - LeNet
 - SimpleMLP
 - Predict
 - DL4J Feedforward Predictor (Classification)
 - DL4J Feedforward Predictor (Layer)
 - DL4J Feedforward Predictor (Regression)
 - Word Embeddings
 - I/O
 - Word Vector Model Reader
 - Word Vector Model Writer
 - Doc2Vec Learner
 - Vocabulary Extractor
 - Word Vector Apply
 - Word2Vec Learner
 - Keras
 - DL Keras Network Learner
 - DL Keras Network Reader
 - DL Network Executor
 - DL Python Network Creator



Keras

TensorFlow™

DEEPLARNING4J

KNIME Extensions

For discussions related to KNIME Extensions

■ Text Processing ■ Scripting ■ Reporting ■ Image Processing
■ REST ■ BigData ■ Deep Learning

Community Extensions

For discussions related to extensions developed by the KNIME community

■ RDKit ■ Scripting Extensions ■ Indigo ■ Erlwood ■ HCS Tools
■ Palladian & Selenium ■ CDK ■ OpenMS ■ Vernalis ■ Seqan

Partner Extensions

For discussions related to extensions developed by our KNIME partners

■ JChem Extensions ■ [Schrödinger Extensions](#)

KNIME Server

For discussions related to KNIME Server

Special Interest Groups

For discussions related to various special interest groups

■ Cheminformatics ■ Bioinformatics

KNIME Development

For discussions related to KNIME development

E agora... Pra onde vou!!

- ▶ Aprofundar nos conceitos matemáticos dos algoritmos.
- ▶ Aprofundar nas outras aplicações e nos outros nós disponíveis no KNIME.
- ▶ E...

Se for do seu interesse aprenda também a programar:



+ SQL + BI

+



ou



ou

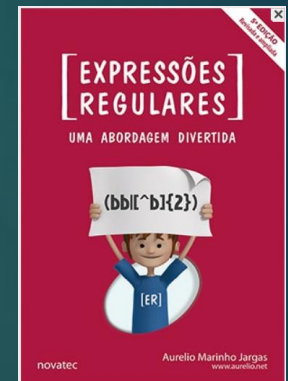
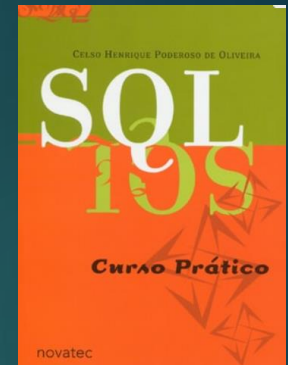
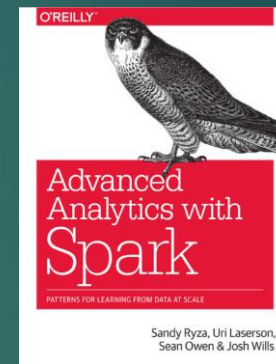
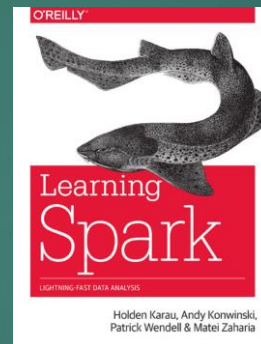
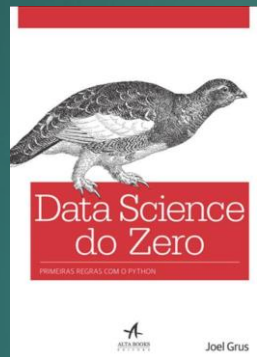
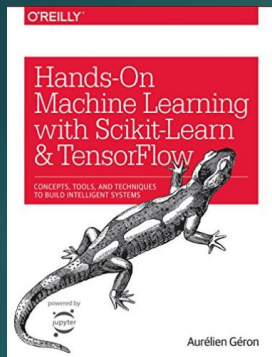
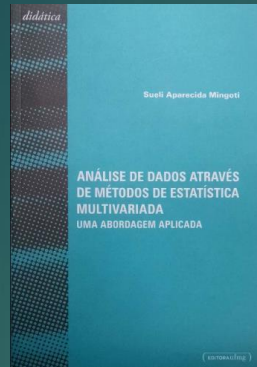
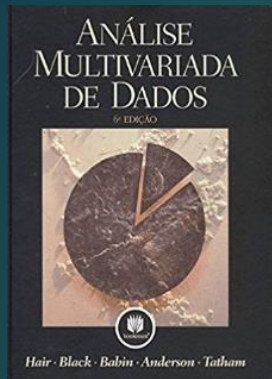


e



e





<http://professores.dcc.ufla.br/~lacerda/>
 Fonte: Prof. Wilian Soares Lacerda
 Depto. Ciência da Computação

- Deep Learning in Python: Master Data Science and Machine Learning with Modern Neural Networks written in Python, Theano, and TensorFlow (Machine Learning in Python) (English Edition) eBook Kindle
- KACHIGAN, Sam Kash. Statistical Analysis: An Interdisciplinary Introduction to Univariate & Multivariate Methods. New York: Radius Press, 1986, 589 p.

“Ninguém sabe tanto que não possa aprender e nem tão pouco que não possa ensinar.”

Anônimo



Gleber Teixeira, PhD
Technical Consultant at Petrobras