

Ciência de Dados (Big Data Processing and Analytics)

Big Data Analytics – Mineração e Análise de Dados



Professor curador
Prof. Dr. Rogério de Oliveira





TRILHA 1

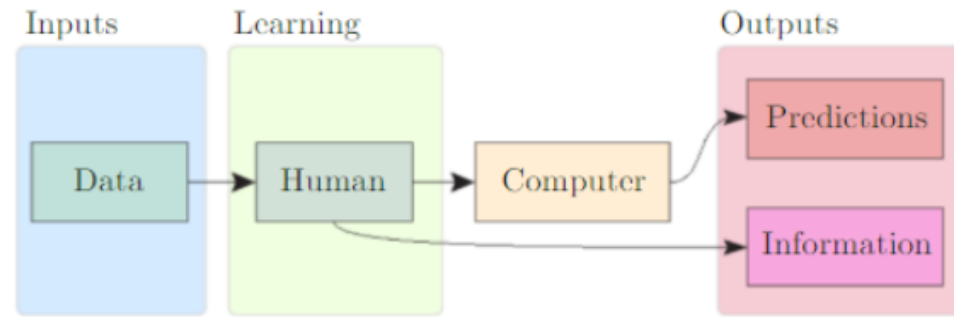
Introdução: Mineração, Ciência de Dados e o Aprendizado de Máquina

Parte A

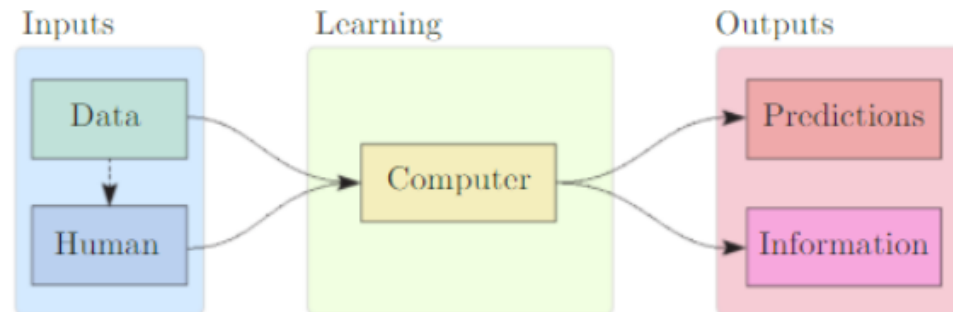
Ciência de Dados e Análise de Dados Tradicional

Tipo	Ciência de Dados	Análise Tradicional
Tipo de Dados	Não estruturados	Estruturados (linhas e colunas)
Volume de Dados	Big Data (centenas de Terabytes)	Dezenas de Terabytes ou menos
Fluxo de Dados	Contínuo, Big Data	Estático
Métodos de Análise	Machine Learning	Visual, Hypothesis-Based
Propósito	Data-base Products	Internal Support Decision

Aprendizado de Máquina – Um novo Paradigma



(a) Without machine learning



(b) With machine learning

Tipos de Padrões e Aplicações

Tarefa	Exemplos
Classificação	Fraud/not Fraud, Churn/Not Churn
Regressão	Preços de Imóveis, Aluguéis, Veículos
Clusterização	Segmentação de Clientes, Produtos, Documentos
Regras de Associação	Pacotes de produtos
Anomaly Detection	Mal funcionamento de dispositivos, Ataques cibernéticos
Matching	Recomendação de produtos ou amigos de uma rede social

Figure 3. Extracting interesting patterns in health outcomes from health-care system use.

Patient	Age	#Medications	Complication
1	52	7	Yes
2	57	9	Yes
3	43	6	Yes
4	33	6	No
5	35	8	No
6	49	8	Yes
7	58	4	No
8	62	3	No
9	48	0	No
10	37	6	Yes



Age ≥ 37
AND
#Medications ≥ 6
→
Complication = Yes (100% confidence)

Exemplos



Machine Learning and Data Science Applications in Industry | <https://github.com/ashishpatel26/Real-time-ML-Project>



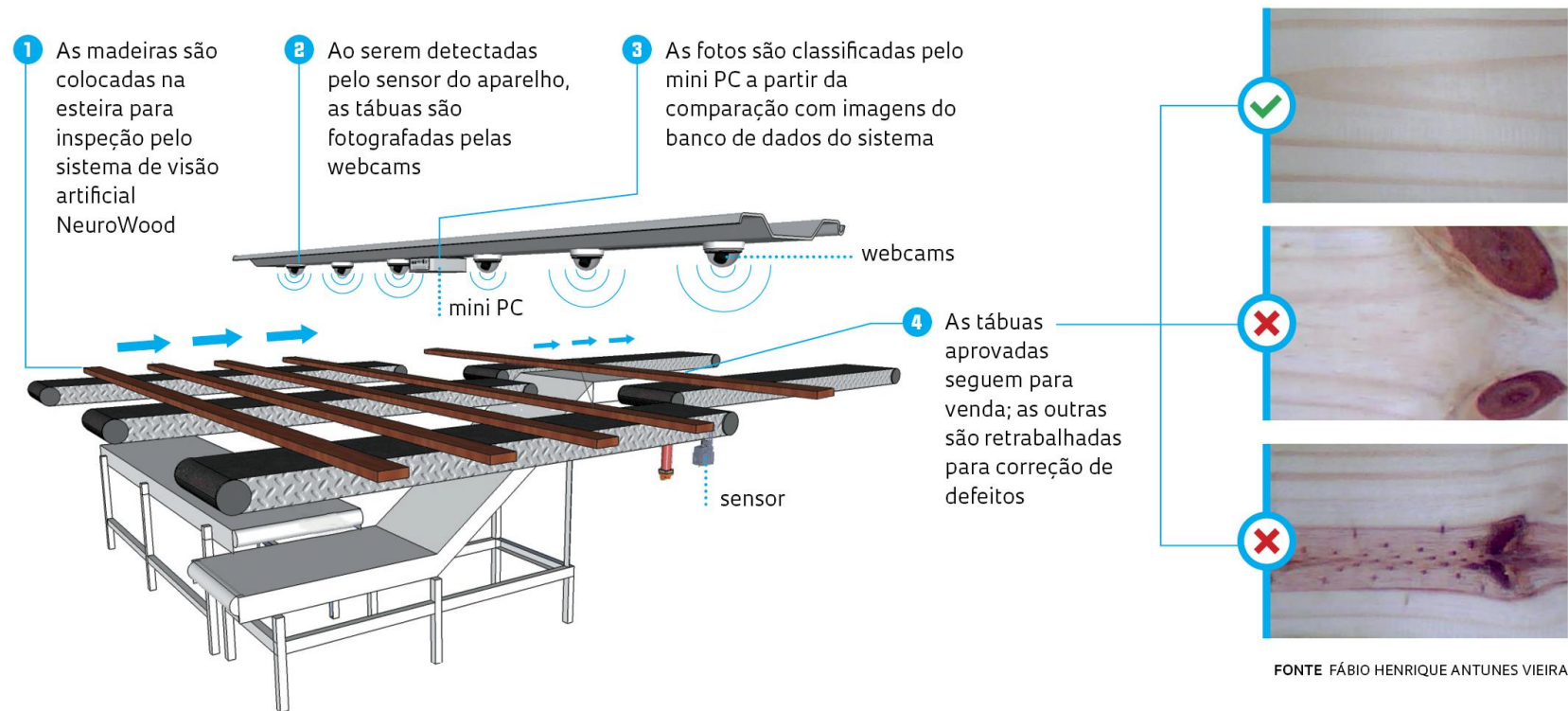
Walmart

Accommodation & Food	Agriculture	Banking & Insurance
Biotechnological & Life Sciences	Construction & Engineering	Education & Research
Emergency & Relief	Finance	Manufacturing
Government and Public Works	Healthcare	Media & Publishing
Justice, Law and Regulations	Miscellaneous	Accounting
Real Estate, Rental & Leasing	Utilities	Wholesale & Retail

Exemplos

De olho na madeira

Saiba como funciona o aparelho que analisa e classifica as tábuas conforme a qualidade



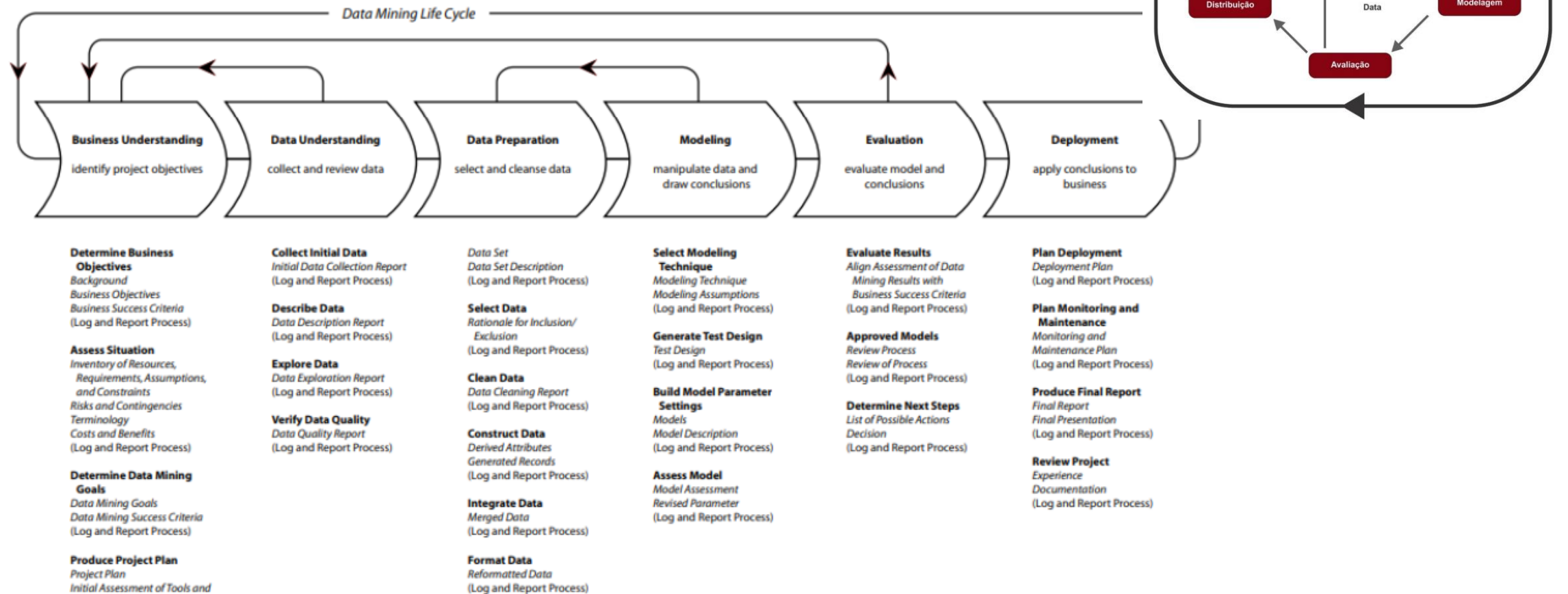


TRILHA 1

Introdução: Mineração, Ciência de Dados e o Aprendizado de Máquina

Parte B

Fases do CRISP-DM



Aprendizado Supervisionado

Tarefas de Aprendizado Supervisionado Breast Cancer Data

Classificação

Árvores de Decisão
Regressão Logística
K-Vizinhos mais Próximos
Support Vector Machines

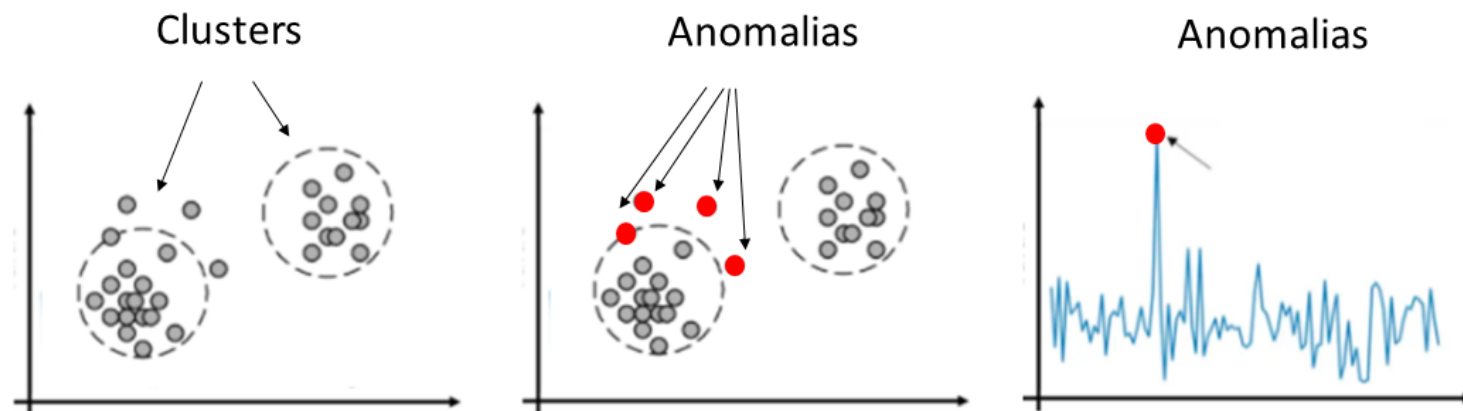
diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean
M	15.340	14.26	102.50	704.4
B	12.880	28.92	82.50	514.3
M	17.080	27.15	111.20	930.9
B	16.140	14.86	104.30	800.0
M	13.480	20.82	88.40	559.2
B	14.470	24.99	95.81	656.4
B	12.490	16.85	79.19	481.6
M	23.210	24.97	153.50	1670.0
B	11.620	18.18	76.38	408.8
B	9.787	19.94	62.11	294.5
M	21.750	20.99	147.30	1491.0
B	10.800	21.98	68.79	359.9
M	25.730	17.46	174.20	2010.0
B	11.870	21.54	76.83	432.0
B	7.691	25.44	48.34	170.4

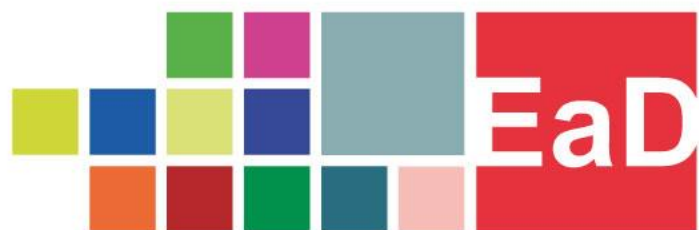
Regressão

Regressão Linear
Regressão Polinomial
Modelos Neurais para Regressão

Aprendizado Não Supervisionado

Exemplos de Aprendizado Não Supervisionado





Universidade Presbiteriana
Mackenzie