

# Machine Learning no Mundo Imperfeito

Marina Fouto



**Mike Townsend**

@Mikettownsend



Follow

Machine learning is like highschool sex. Everyone says they do it, nobody really does, and no one knows what it actually is.

2:04 PM - 26 Sep 2016



376



452

# O que é Machine Learning?

Algoritmos capazes de “aprender” a partir de dados empíricos e fazer previsões sobre estes

# Conceitos Básicos I

Classes

Features (Características)

Datasets

# Conceitos Básicos II

## Aprendizado Supervisionado

Classificação: divide os dados em classes, precisa de treinamento

Regressão: infere novos dados a partir de observações passadas

## Aprendizado Não-Supervisionado

Clusterização: agrupamento de dados similares

Aproximação: “pessoas que compraram X também compraram Y”

# Machine Learning... Por Onde Começar?

Qual é o problema que você quer resolver?

Você sabe o que são os seus dados?

# Conseguindo Dados I

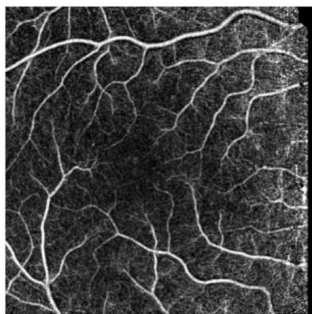
Seus dados podem ser:

Organizados e lindos como o Iris dataset

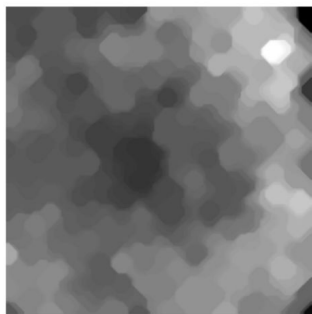
Cheios de observações faltantes

Você pode ter que extrair eles de algo

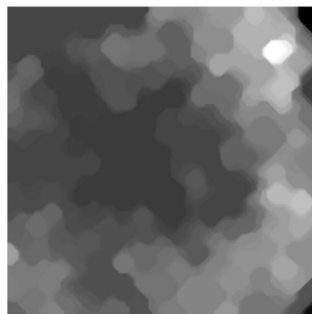
# Conseguindo Dados II



(a)



(b)



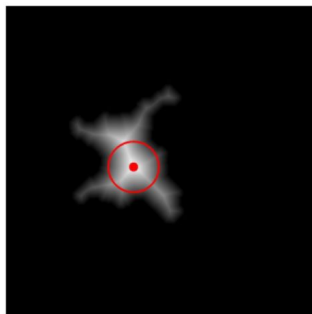
(c)



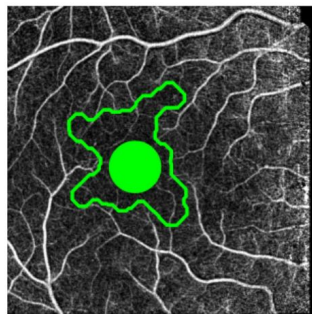
(d)



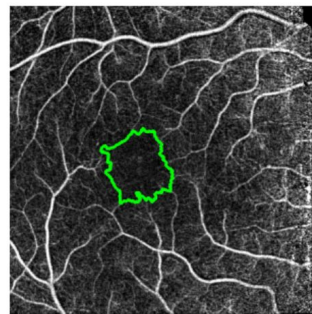
(e)



(f)



(g)

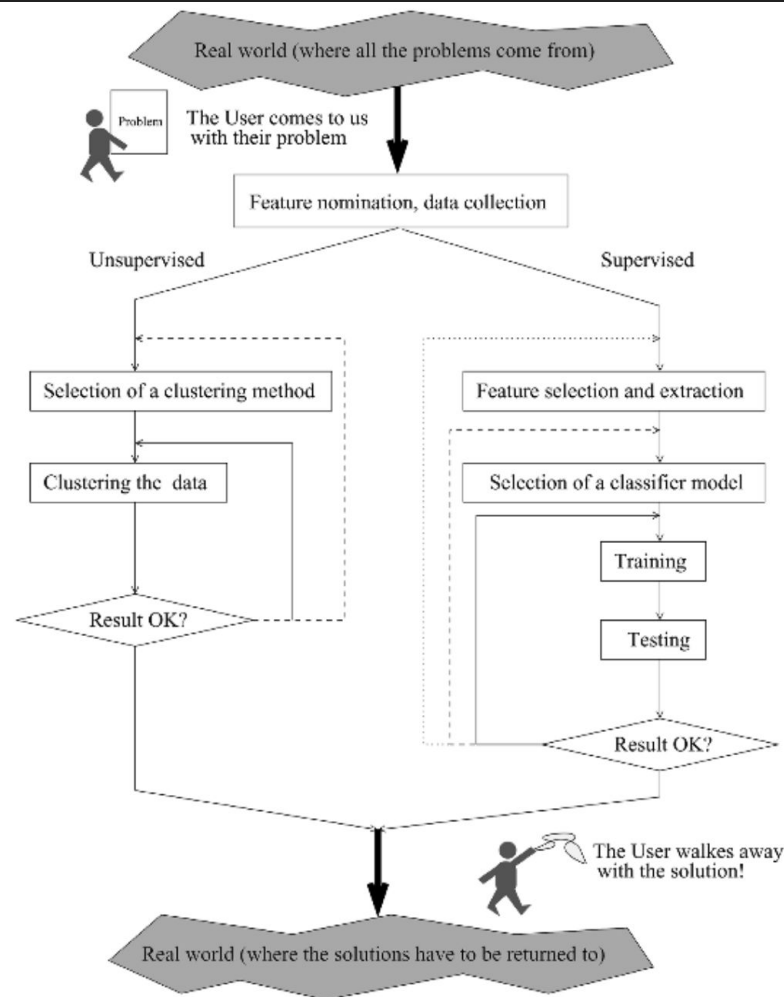


(h)

19 características extraídas baseadas no formato da região obtida pela segmentação

7 datasets construídos a partir da combinação das características





Kuncheva, 2004

# Aprendizado Supervisionado - Classificação

Support Vector Machines (SVM)

K-Nearest Neighbors (KNN)

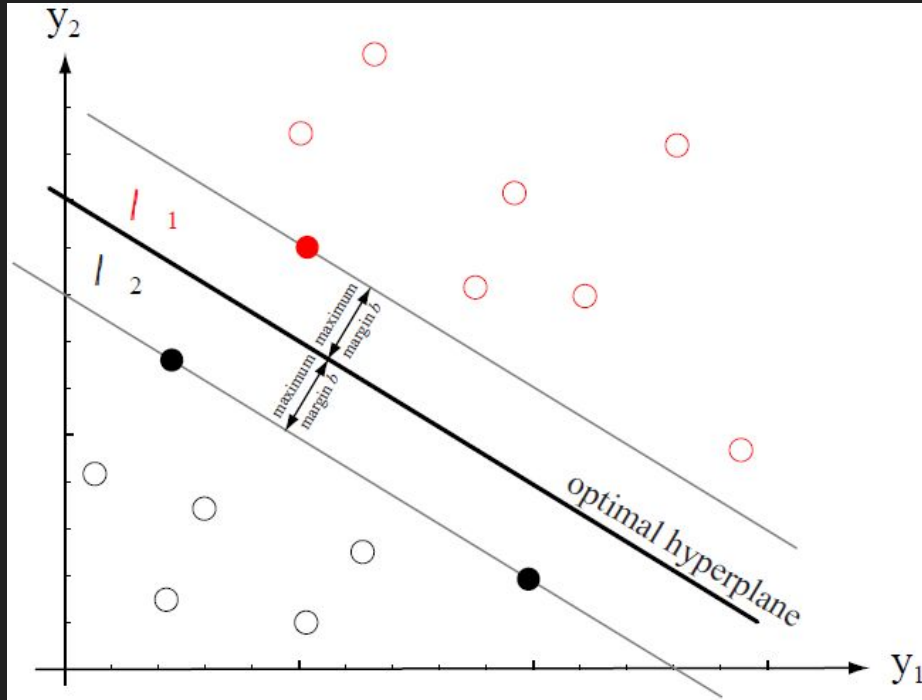
Naive Bayes

Árvores de Decisão

Ensemble Methods

Redes Neurais

# Support Vector Machines (SVM)

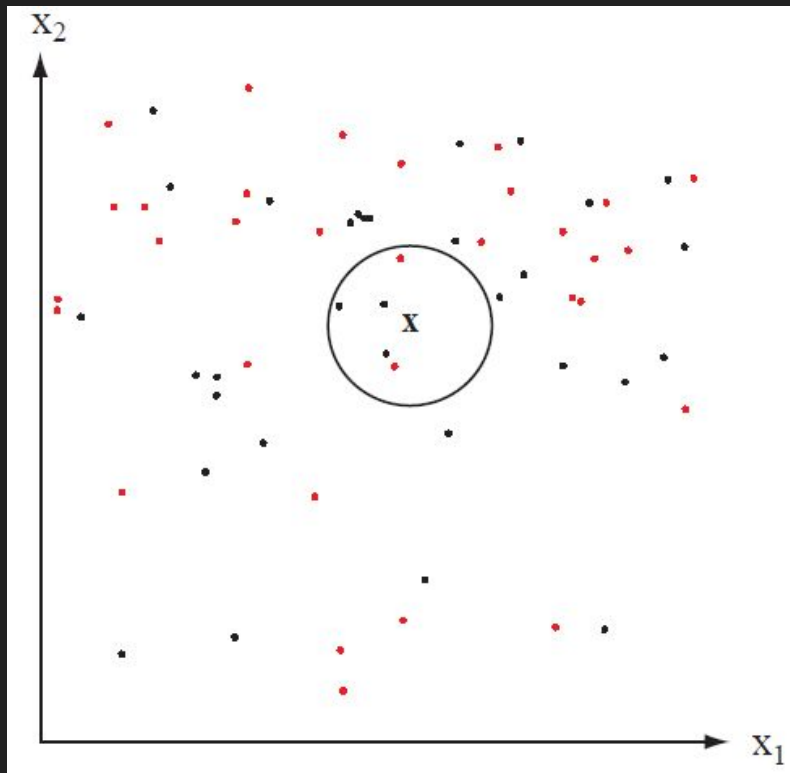


Classificador binário

Versátil, pode utilizar várias funções como separador

Útil em espaços de alta dimensionalidade e eficaz na classificação quando a quantidade de dimensões é maior que a de amostras

# K-Nearest Neighbors (KNN)



KNN classifica a amostra a partir de voto majoritário

Escolhe-se vizinhança ímpar para evitar empates

# Naive Bayes

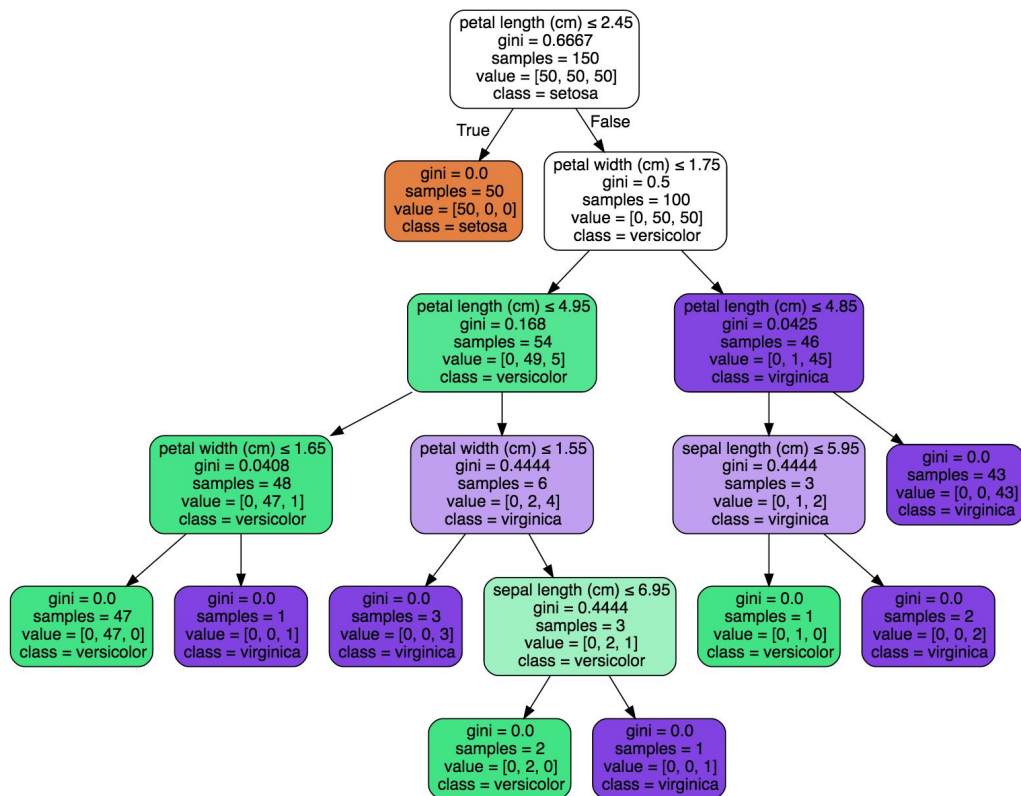
Naive = Ingênuo

Assume que os atributos são condicionalmente independentes em uma mesma classe

Muito utilizado em detecção de spam e é um bom classificador quando existem poucos dados a serem analisados

(<http://www.cs.unb.ca/~hzhang/publications/FLAIRS04ZhangH.pdf>)

# Árvores de Decisão



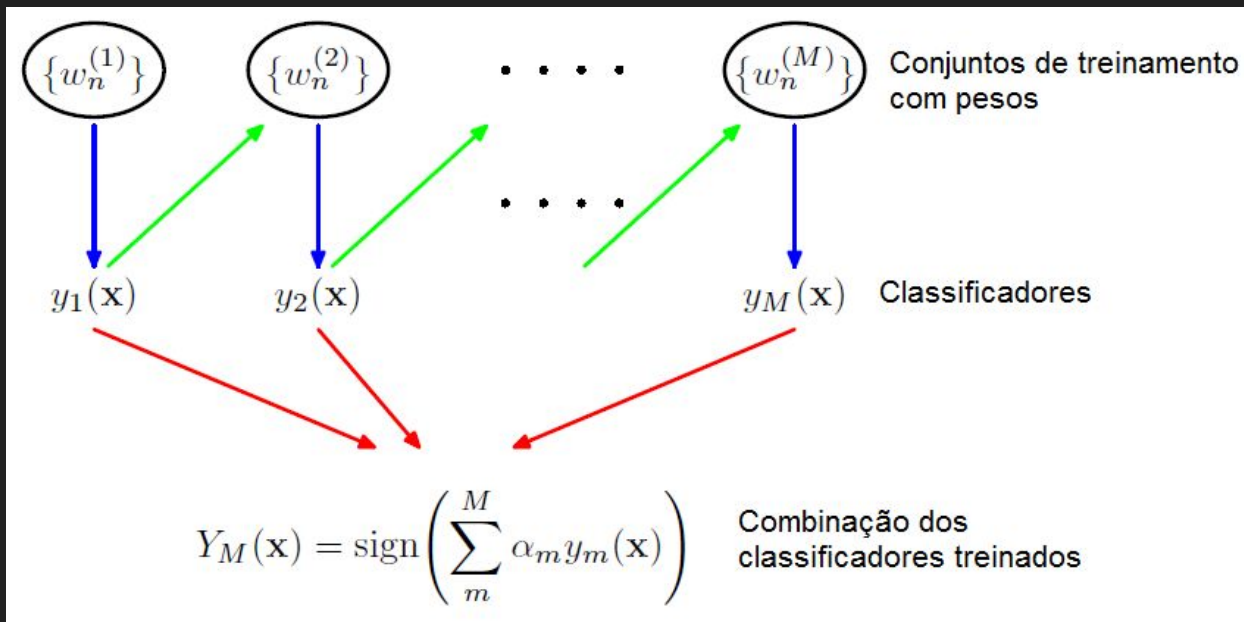
Fáceis de entender e visualizar

Multi-classe

Podem gerar árvores muito complexas e levar a overfitting

Pode favorecer uma classe, caso ela seja maioria no dataset

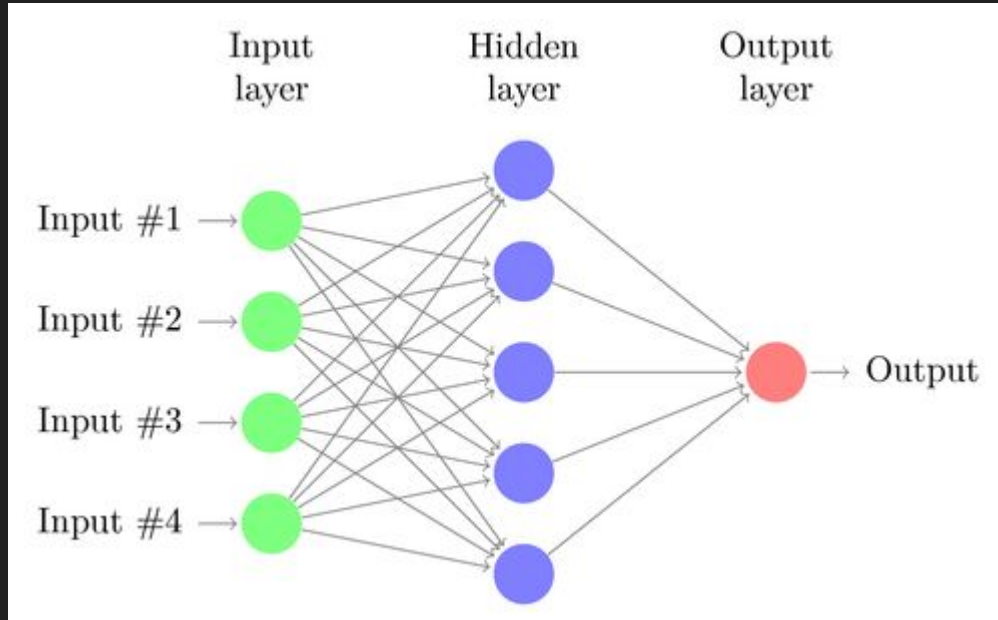
# Ensemble Methods



Métodos de média: a média dos classificadores aplicados independentemente sobre os dados é tomada como resultado

Métodos de boosting: um classificador é construído após o outro, baseado no resultado anterior

# Redes Neurais



Multi-layer Perceptron (MLP)

Redes Neurais Convolucionais (CNN)

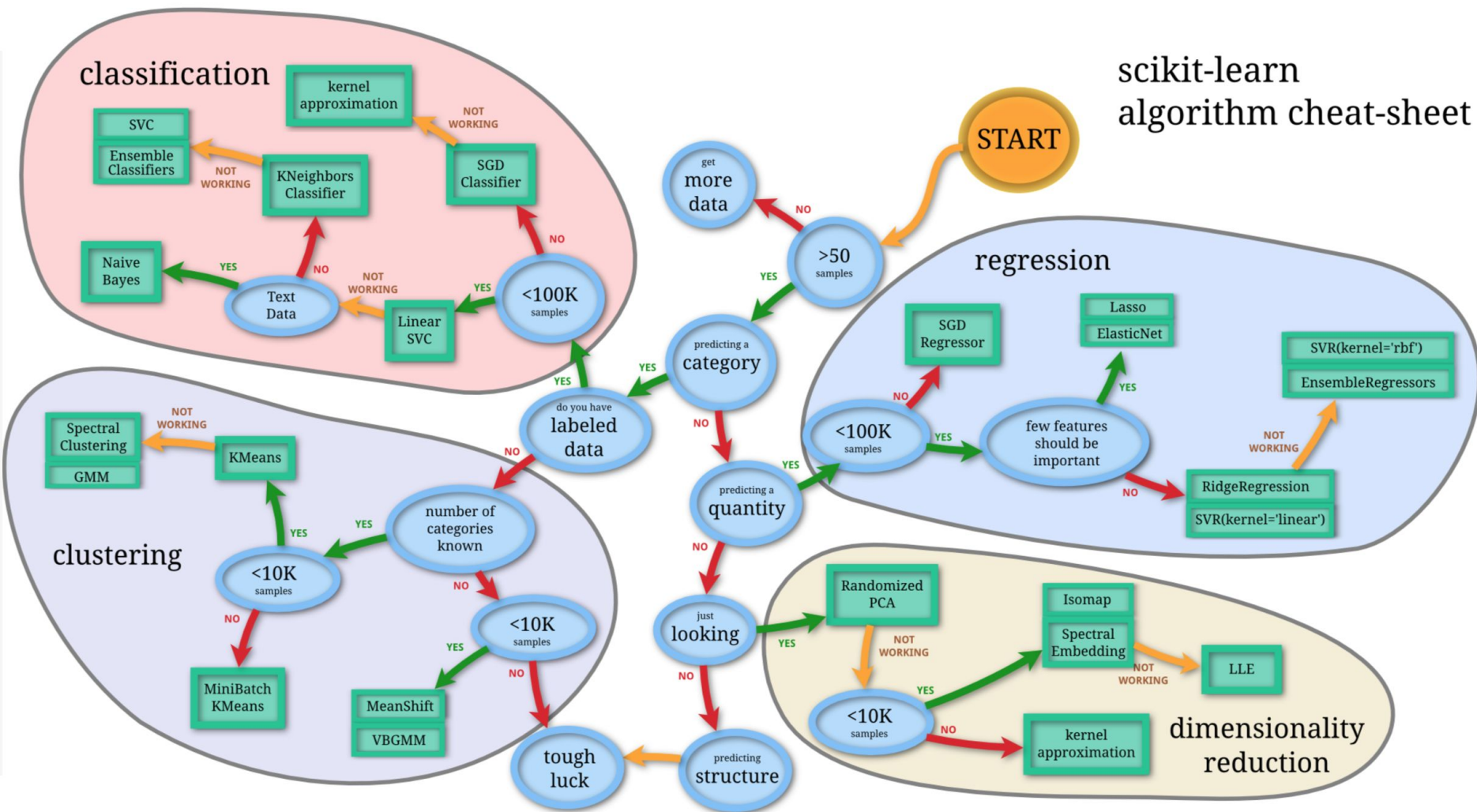
Redes Neurais Recorrentes (RNN)

Redes Neurais Probabilísticas (PNN)



# Aprendizado Não Supervisionado - Clusterização

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
<a href="#">K-Means</a>	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> with <a href="#">MiniBatch code</a>	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
<a href="#">Affinity propagation</a>	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
<a href="#">Mean-shift</a>	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
<a href="#">Spectral clustering</a>	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
<a href="#">Ward hierarchical clustering</a>	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
<a href="#">Agglomerative clustering</a>	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
<a href="#">DBSCAN</a>	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes	Distances between nearest points
<a href="#">Gaussian mixtures</a>	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
<a href="#">Birch</a>	branching factor, threshold, optional global clusterer.	Large <code>n_clusters</code> and <code>n_samples</code>	Large dataset, outlier removal, data reduction.	Euclidean distance between points

scikit-learn  
algorithm cheat-sheet

# Redução de Dimensionalidade

Maldição da dimensionalidade

Serve para datasets com muitas dimensões, onde é difícil entender e avaliar os dados

Principal Component Analysis (PCA), decompõe um dataset de alta dimensionalidade em um conjunto de componentes ortogonais que explicam a maior quantidade de variância

# Dados desbalanceados

Datasets médicos tendem a ter dados desbalanceados

Não confiar somente em % de acerto. 95% de acerto pode não ser útil

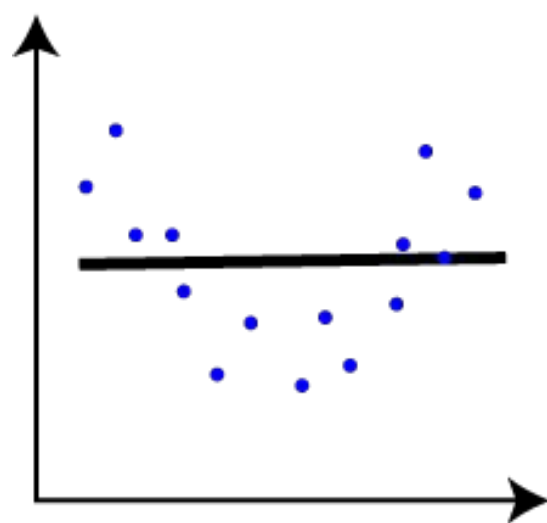
Sensitividade (Verdadeiro Positivo)

Especificidade (Verdadeiro Negativo)

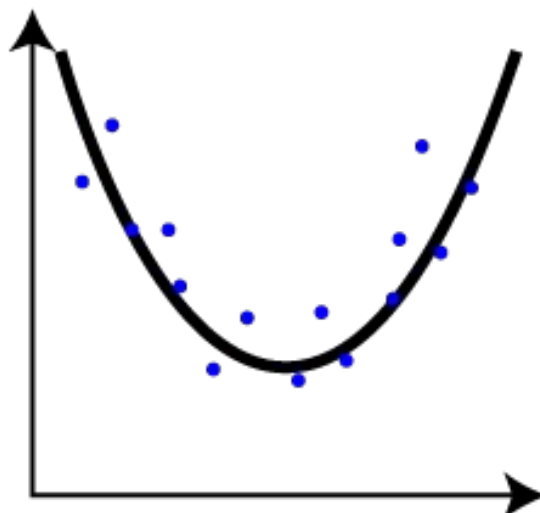
Balancear os dados com oversampling ou undersampling

# Overfitting

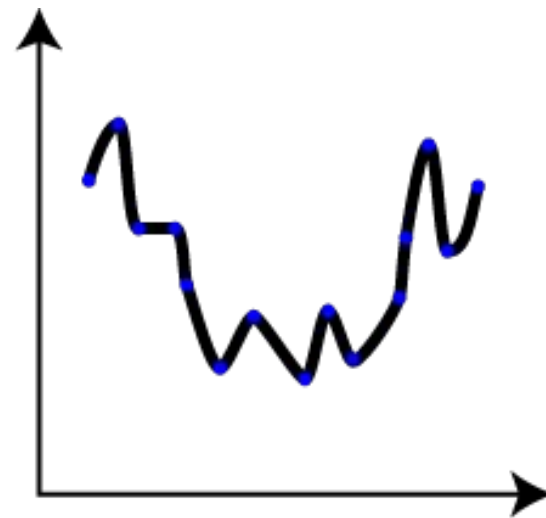
Um modelo flexível demais não generaliza!



Underfitting



Adequado



Overfitting

# Cross Validation

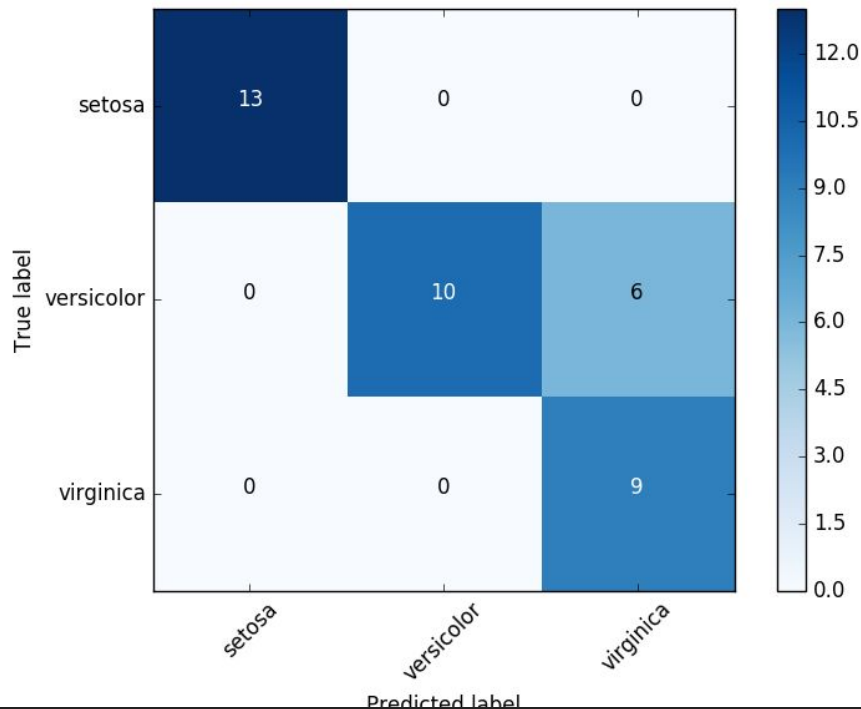
Evita overfitting

Divide o conjunto de treinamento em treinamento e validação

2-fold (50% treinamento, 50% validação)

K-fold (treina com  $k-1$  partes, valida com 1)

# Matriz de Confusão



	P' (Predicted)	N' (Predicted)
P (Actual)	True Positive	False Negative
N (Actual)	False Positive	True Negative

# Type I e Type II

Um erro Type I acontece quando um valor negativo é classificado como positivo

Um erro Type II é quando um valor positivo é classificado como negativo



# TPOT

<https://github.com/rhiever/tpot>

Biblioteca escrita em Python para otimizar pipeline de machine learning utilizando algoritmos genéticos

# Para Saber Mais

Kaggle! <https://www.kaggle.com/kernels>

## Welcome to Kaggle Kernels

The best place to explore data science results and share your own work



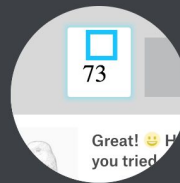
### Code

Skip the download. Kernels is preloaded with the most common data science languages and libraries.



### Learn

Gain exposure to new tools and techniques. The “hottest” kernels showcase the best work on Kaggle.



### Mentor

Give back by sharing what you know. You can answer questions and leave feedback on others' code and results.

# Para Saber Mais

Data Scientist Workbench! <https://datascientistworkbench.com/>

## BUILDING YOUR ANALYTICS

All the analytics tools you need.



**IPYTHON / JUPYTER**

Your choice of Python, R, or Scala notebooks.



**APACHE ZEPPELIN**

Multiple languages in the same notebook.



**RSTUDIO IDE**

A complete statistical suite in the cloud.



**SEAHORSE**

Program Apache Spark visually.

# Para Saber Mais

Introdução a Machine Learning no Big Data University:

<https://bigdatauniversity.com/courses/introduction-to-machine-learning/>

Nanodegree em Machine Learning no Udacity

<https://br.udacity.com/course/machine-learning-engineer-nanodegree--nd009/>

edX: [https://www.edx.org/course?search\\_query=machine+learning](https://www.edx.org/course?search_query=machine+learning)

# Machine Learning no Mundo Imperfeito

Marina Fouto