

Astronomia ao meio-dia

Machine learning e alguns exemplos de aplicação em dados astronômicos

Nina S. T. Hirata



Departamento de Ciência da Computação
Instituto de Matemática e Estatística
Universidade of São Paulo (USP)



São Paulo, 03/05/2018

Máquinas inteligentes?

Existem máquinas inteligentes ?

O que as máquinas inteligentes são capazes de fazer ?

Máquinas aprendem ? (elas podem fazer mais ?)

Máquinas inteligentes?

Existem máquinas inteligentes ?

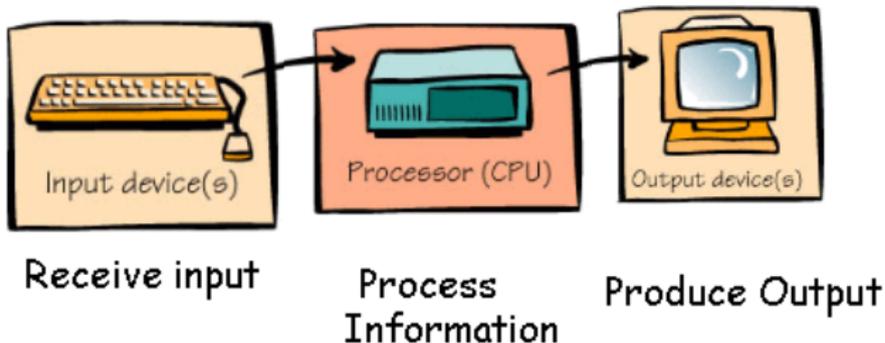
O que as máquinas inteligentes são capazes de fazer ?

Máquinas aprendem ? (elas podem fazer mais ?)

Não vamos discutir este tema do ponto de vista filosófico / ético

Vamos discutir ML do ponto de vista computacional

O que os computadores fazem ?



<http://gebar.weebly.com/inputoutput-devices.html>

Computadores processam dados segundo um **algoritmo**

Algoritmos são soluções para **problemas computacionais**

O que são **problemas computacionais** ?

Exemplo: somar números

- Dada uma lista de números, calcular a soma dos números na lista

Entrada \Rightarrow Saída

3, 1, 7 \Rightarrow 11

0.25, 0.75, 0.5, 0.1 \Rightarrow 1.6

1, 3, 5, 7, 9 \Rightarrow 25

Algoritmo

Entrada: uma lista de números a serem somados

Saída: a soma (total) dos números na lista

SOMA = 0

Enquanto lista de números não está vazia

 Coloque o próximo número da lista em NUM
 (remove o número da lista)

 SOMA = SOMA + NUM

Imprima SOMA

Exemplo: o que está sendo computado neste caso ?

Entrada \Rightarrow Saída

9, 2, 0, -1, 7, 4 \Rightarrow -1, 0, 2, 4, 7, 9

'x', 'a', 'm', 'b' \Rightarrow 'a', 'b', 'm', 'x'

Exemplo: o que está sendo computado neste caso ?

Entrada \Rightarrow Saída

9, 2, 0, -1, 7, 4 \Rightarrow -1, 0, 2, 4, 7, 9

'x', 'a', 'm', 'b' \Rightarrow 'a', 'b', 'm', 'x'

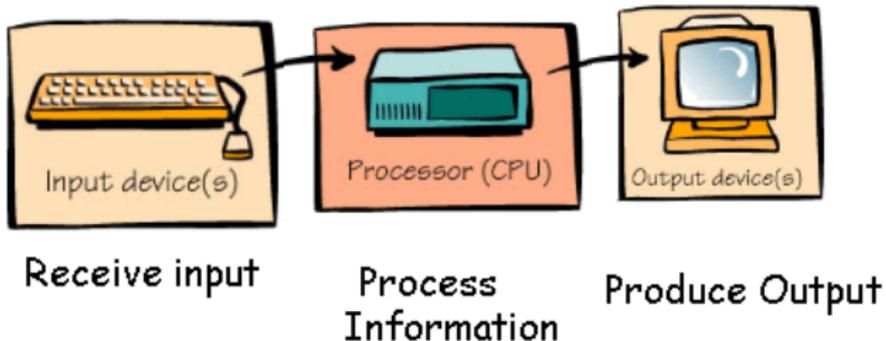
Problema de ordenação!

Nós sabemos como resolver (computar) esse problema!

Vários algoritmos de ordenação

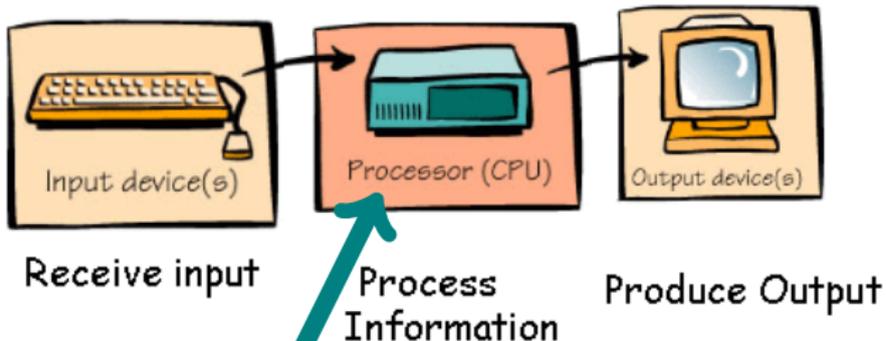
- selection sort
- insertion sort
- bubble sort
- merge sort
- quick sort
- etc

Como os dados são processados?



<http://gebar.weebly.com/inputoutput-devices.html>

Como os dados são processados?



<http://gebr.weebly.com/inputoutput-devices.html>

Programas (implementação de algoritmos)

Desafios computacionais

- Garantir que o algoritmo está correto
- Garantir que a implementação do algoritmo está correta
- Desenhar algoritmos eficientes
- Fazer implementações eficientes

- O que fazer quando não há algoritmo eficiente ?
- O que fazer quando a relação entrada-saída não está bem definida?
- Ou quando não sabemos descrever formalmente a relação?

Quando não temos um algoritmo

 \Rightarrow 'a'

 \Rightarrow 'A'

 \Rightarrow 'b'

 \Rightarrow 'X'

 \Rightarrow 'd'

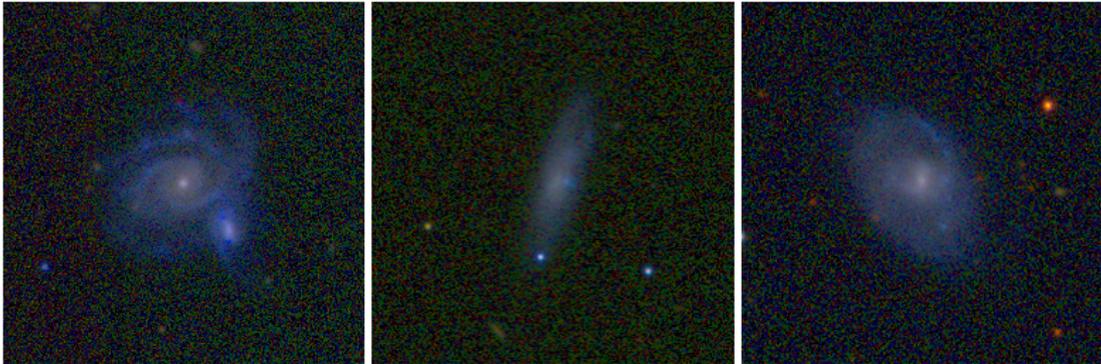
 \Rightarrow 'M'

Quando não temos um algoritmo

Elíptica

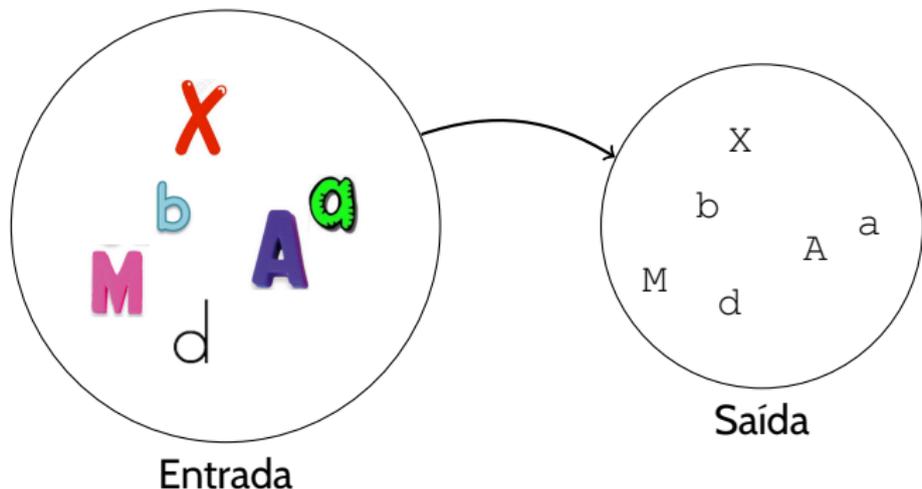


Espiral

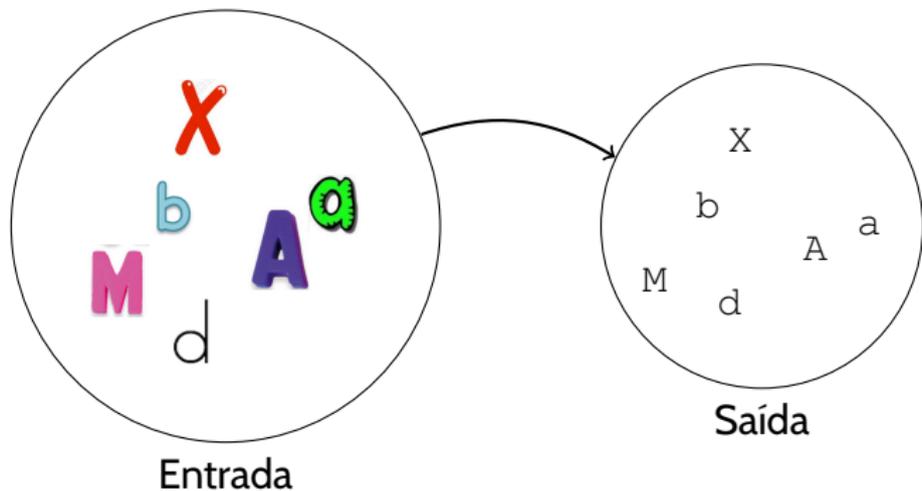


Imagens do EFIGI

Entra em cena ML



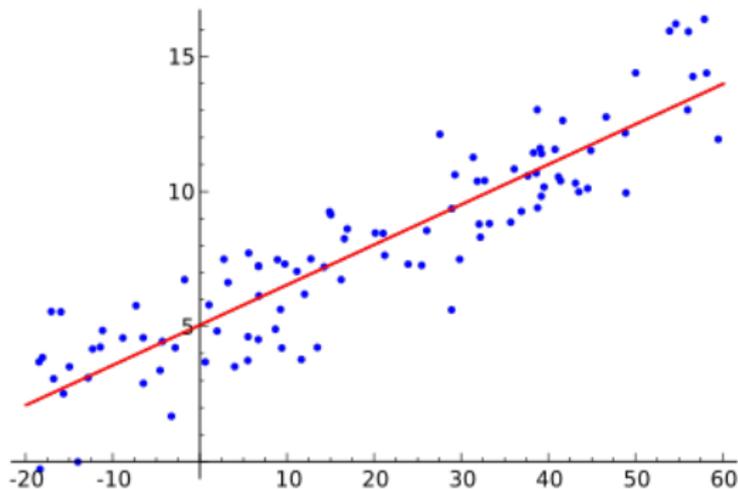
Há uma **relação** entre elementos no espaço de entrada e elementos no espaço de saída \implies que eu não sei descrever formalmente



Idéia central de ML: “aprender” a relação a partir de exemplos

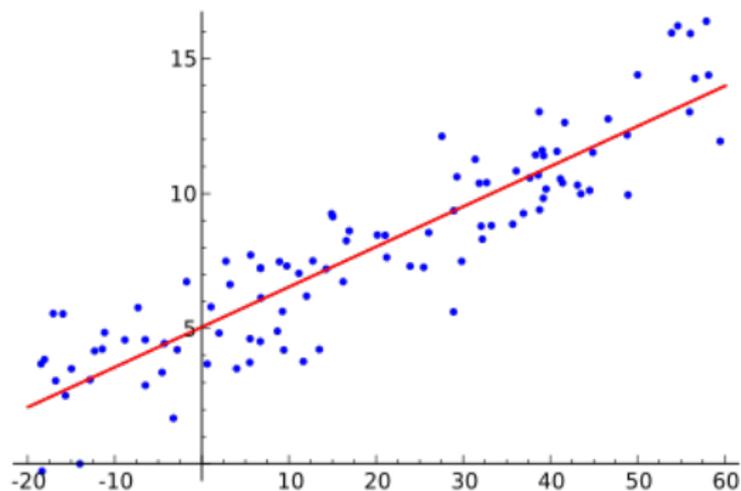
Produto final de ML: um programa que realiza o processamento entrada-saída

Um exemplo simples: regressão linear



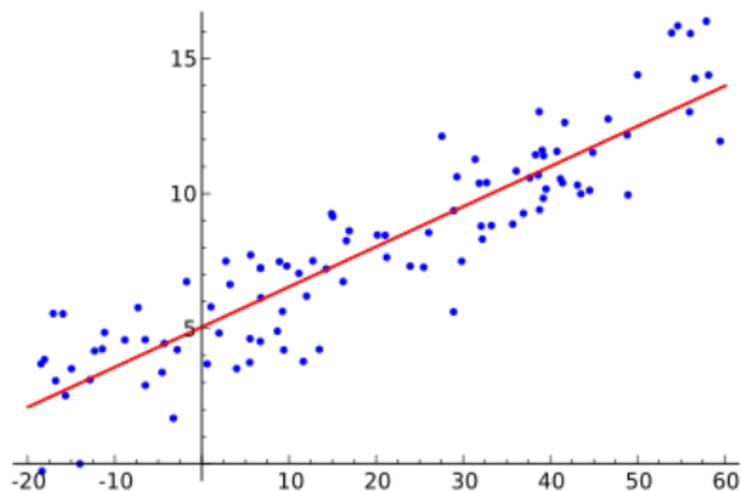
Os pontos azuis (x_i, y_i) são os exemplos de treinamento

Um exemplo simples: regressão linear



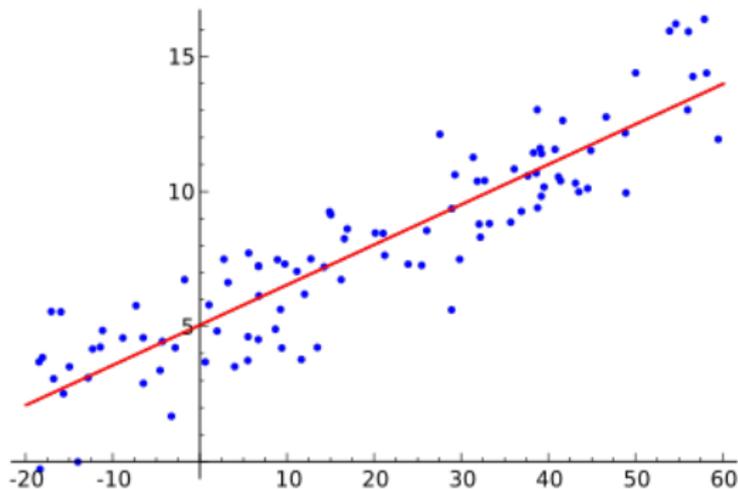
Há uma **relação linear** entre x e y

Um exemplo simples: regressão linear



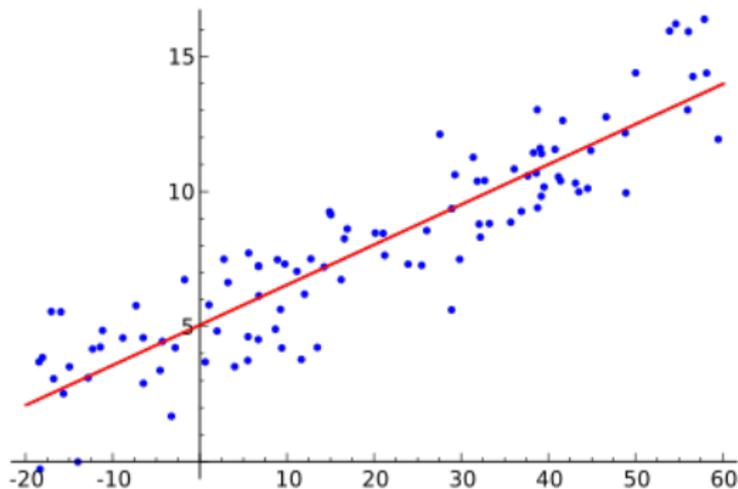
Família de hipóteses adequada: $g(x) = a + b x$

Um exemplo simples: regressão linear



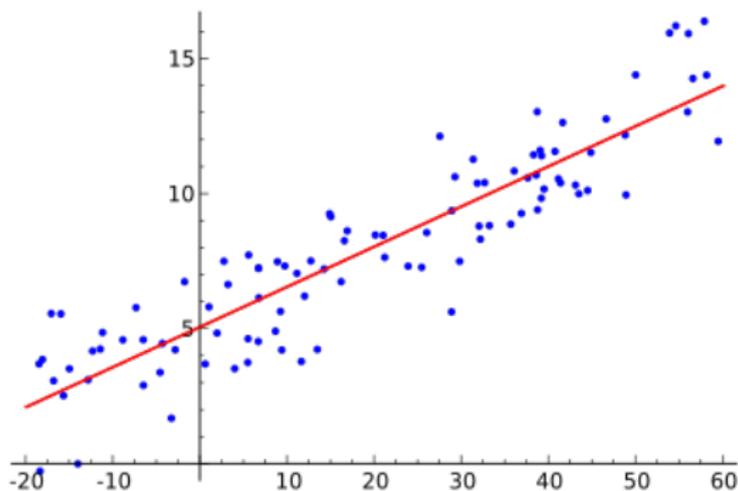
Função custo : $J(a, b) = \frac{1}{N} \sum_{i=1}^N (y_i - g(x_i))^2$

Um exemplo simples: regressão linear



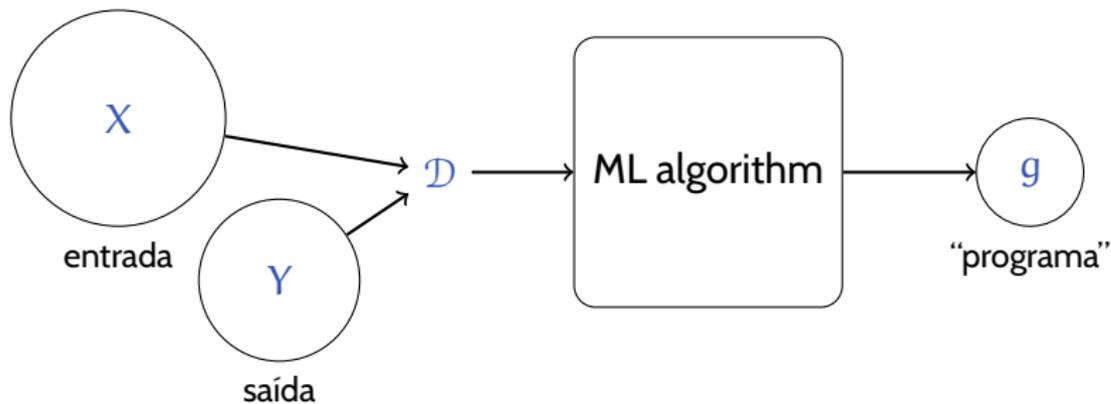
Otimização : encontrar a e b ($g(x) = a + b x$) que minimiza o custo $J(a, b)$

Um exemplo simples: regressão linear

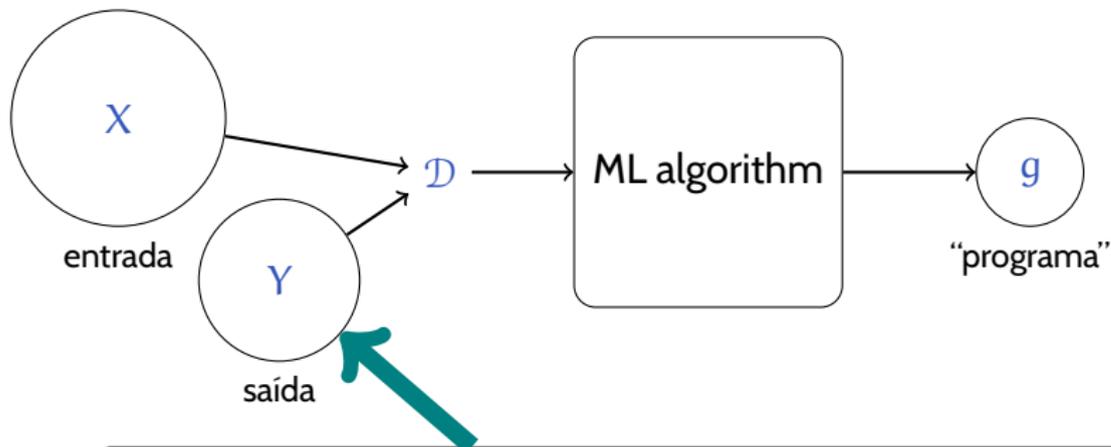


Aprendizado : há uma solução analítica para a regressão linear mas pode-se também usar uma abordagem de ajuste de parâmetros iterativa (ex.: *gradient descent*)

Componentes de um processo de ML

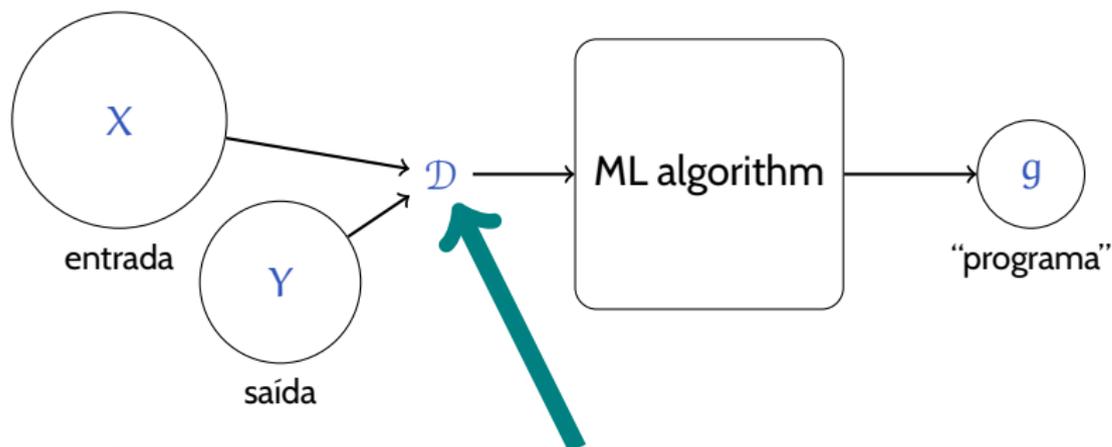


Componentes de um processo de ML



Entrada-saída caracteriza o processamento desejado

Componentes de um processo de ML

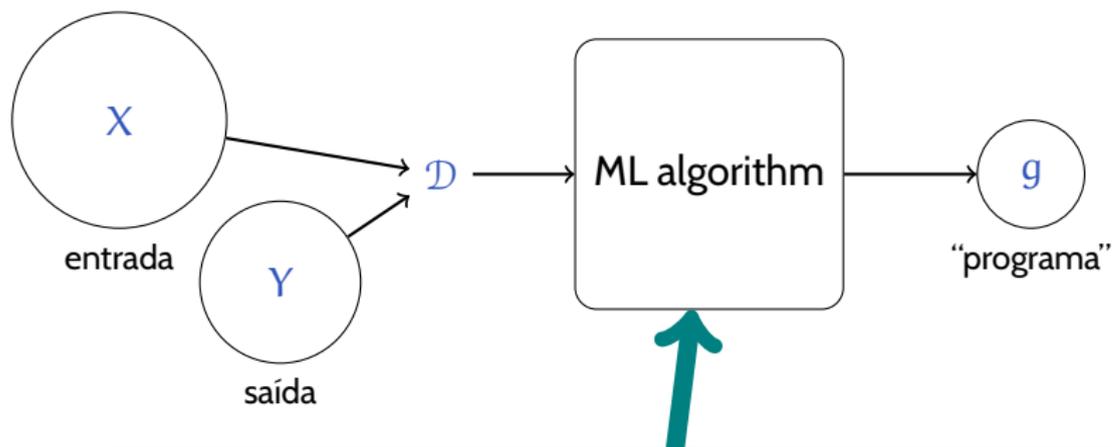


Dados de treinamento

$$\mathcal{D} = \left\{ (\mathbf{x}_i, y_i) \in X \times Y : i = 1, 2, \dots, N \right\}$$

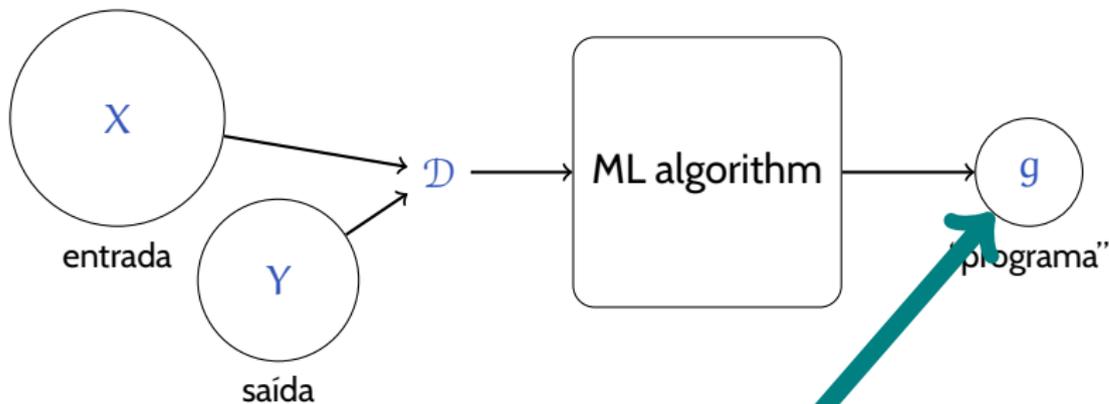
(exemplos de *relações input-output*)

Componentes de um processo de ML



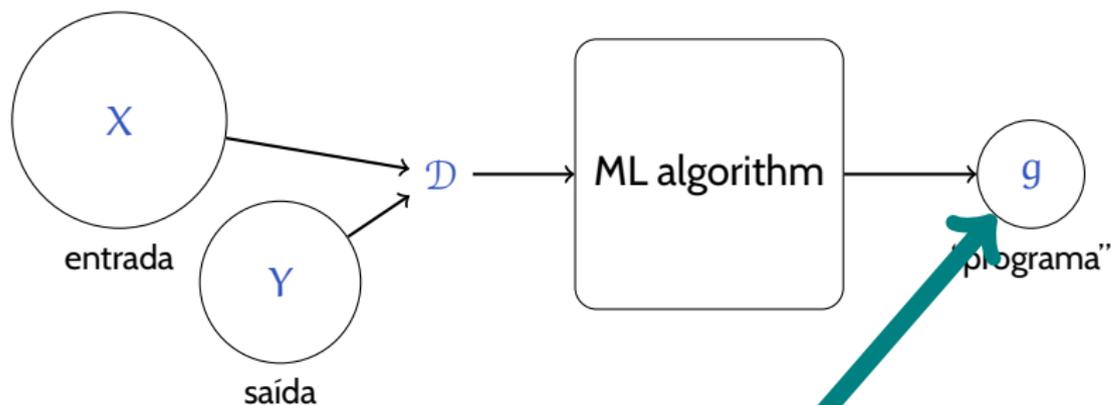
Algoritmo de aprendizado usa exemplos em \mathcal{D}
para produzir um “programa” g

Componentes de um processo de ML



O “programa”, ou hipótese $g : X \rightarrow Y$, é o que queremos
Deve ser tal que $\hat{y} = g(\mathbf{x})$ seja
“o mais próximo possível” de y

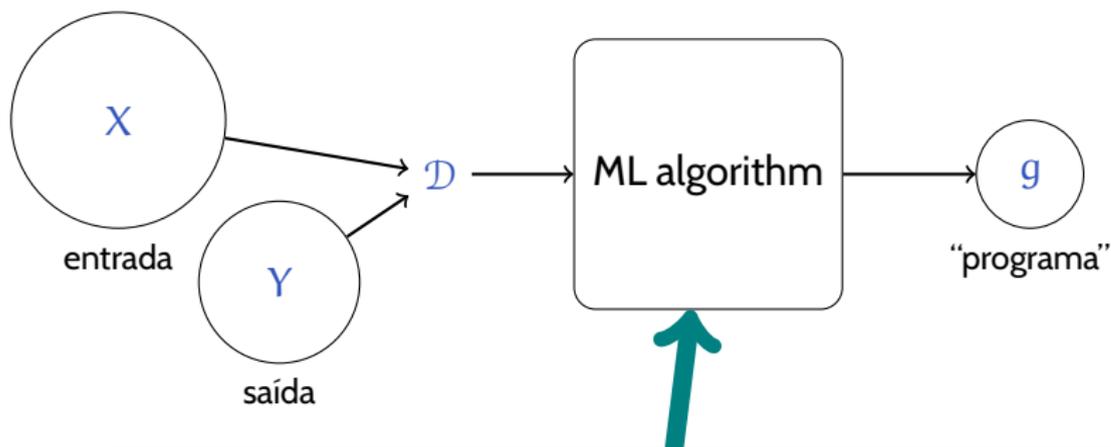
Componentes de um processo de ML



Na prática, tenta-se encontrar $\hat{y} = g(\mathbf{x})$ que minimiza um erro empírico; ex.:

$$\hat{E}_{\text{err}}(f) = \frac{1}{N} \sum_{i=1}^N (y_i - g(\mathbf{x}_i))^2$$

Componentes de um processo de ML



Portanto, em geral os algoritmos de aprendizado usam alguma **técnica de otimização**

Exemplo

Classificação de imagens de galáxia

O que precisamos considerar

Entrada / saída

X : todas as possíveis imagens de galáxias (espiral / elíptica)

Y : classe (espiral / elíptica)

Dados de treinamento

Quantos pares (\mathbf{x}, y) precisamos? De onde eles vem ?

Treinamento

Custo a ser minimizado : erro de classificação

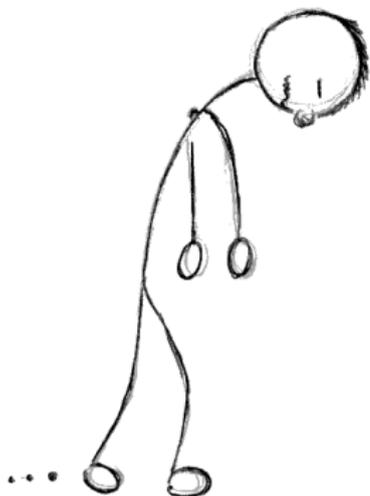
Algoritmo de aprendizado : SVM, redes neurais, random forests, etc

Representação dos dados

Como esses algoritmos irão “enxergar” as imagens ? \implies como representá-las?

Representação dos dados de entrada

Exemplo: No caso de diagnóstico médico, informações sobre o paciente

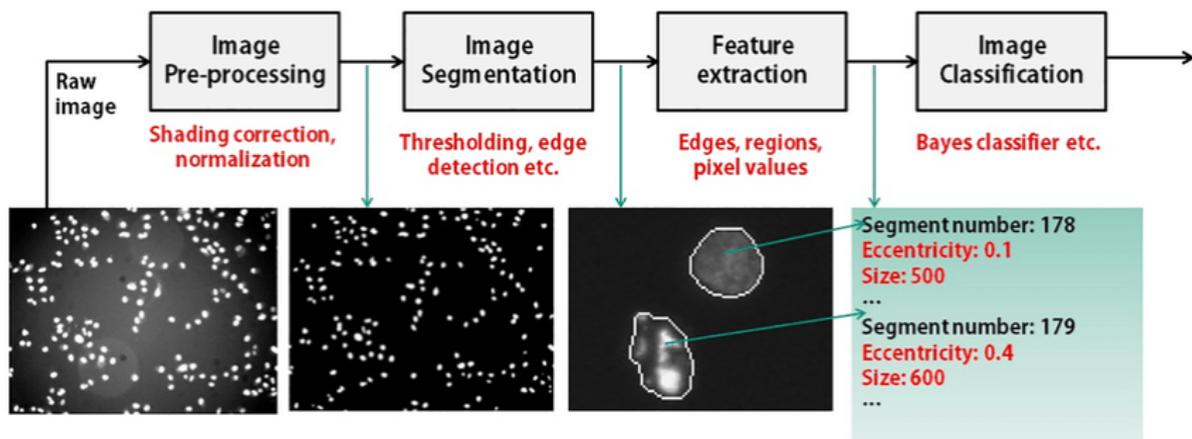


- x_1 pressão arterial
- x_2 temperatura
- x_3 sente dor de cabeça
- x_4 peso
- x_5 altura
- x_6 idade
- x_7 sente fadiga
- x_8 tem tosse
- x_j ...
- x_n etc

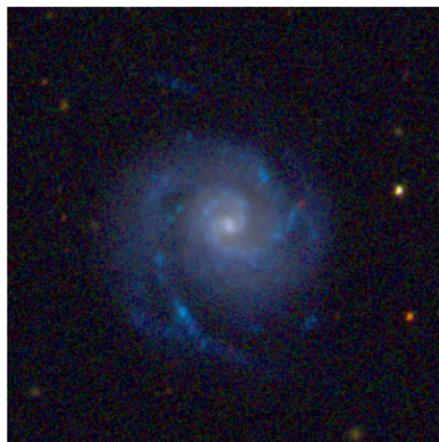


$$\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$$

Quando a entrada é uma imagem



Classificação de imagens de galáxias



Processamento de imagens

Processamento de sinais

Extração de características diversas

- concentration
- asymmetry
- smoothness
- entropy
- spirality, etc

$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

Classificação de imagens de galáxias

(com Fabrício Ferrari e Mateus Espadoto)

Dataset usado: EFIGI (<https://www.astromatic.net/projects/efigi>)

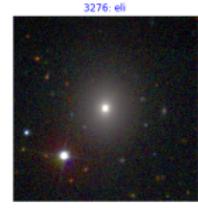
Procedimento executado:

1. Selecionar apenas as espirais e as elípticas
2. Extrair *features* para cada imagem
 - Usando Morfometryka (dados fornecidos por Fabrício Ferrari, FURG)
 - Usando uma rede convolucional pré-treinada com o ImageNet
3. Separar dados em treinamento e validação
4. Treinar SVM usando dados de treinamento
5. Avaliar desempenho (de predição) no conjunto de validação

Features convolucionais

Classificação errada

4 vizinhos mais próximos



Features convolucionais

Classificação errada

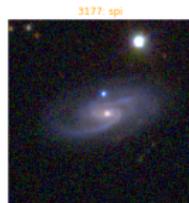
4 vizinhos mais próximos



Features morfológicos

Classificação errada

4 vizinhos mais próximos



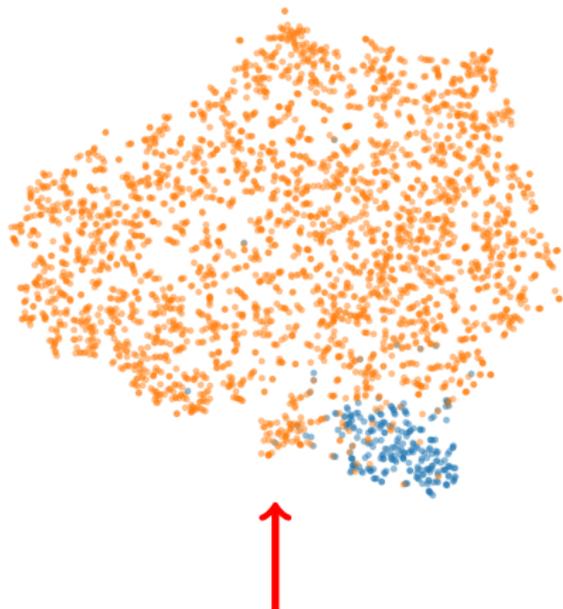
Features morfológicos

Classificação errada

4 vizinhos mais próximos

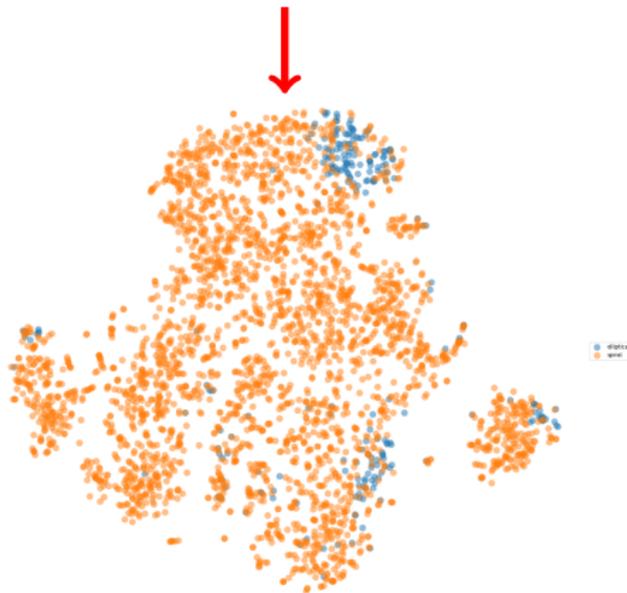


Projeções das features

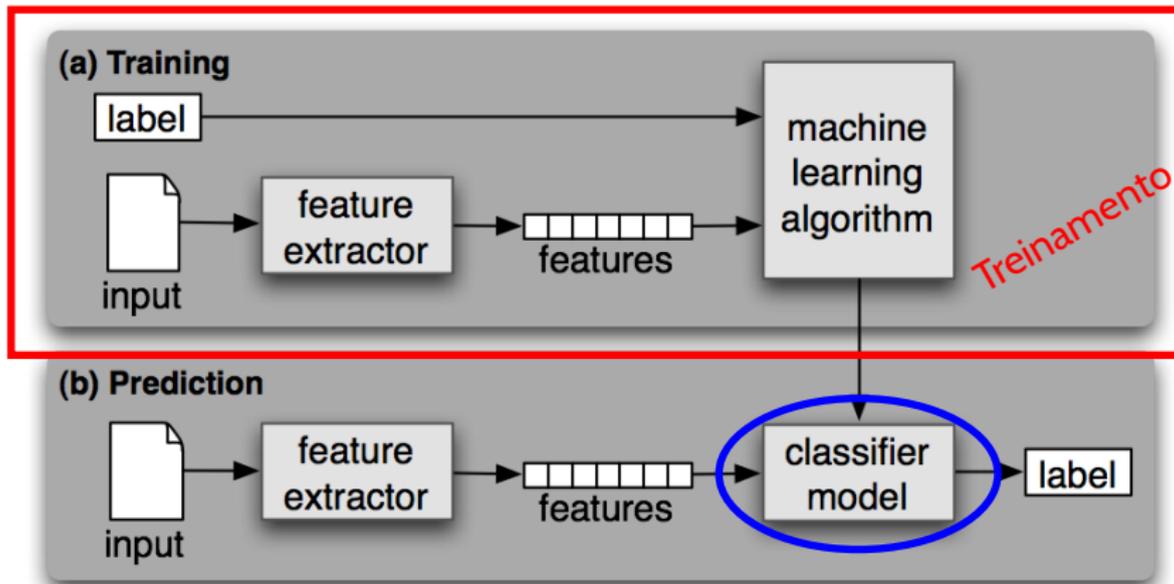


Features convolucionais

Features morfométricos

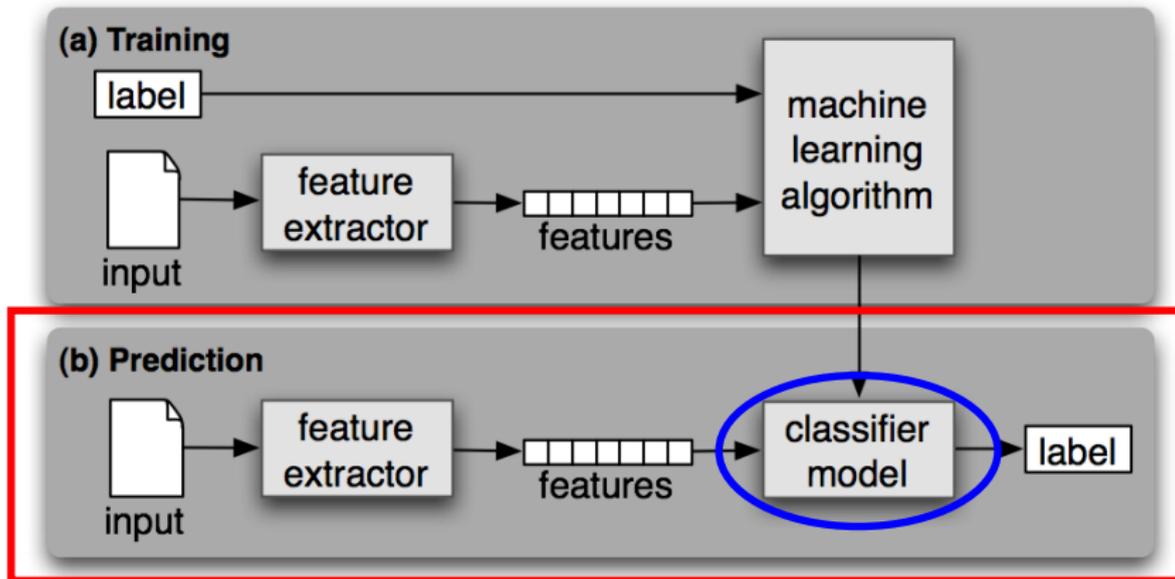


Processo de machine learning



<http://www.nltk.org/book/ch06.html>

Processo de machine learning



<http://www.nltk.org/book/ch06.html>

Predição

Alguns comentários

- ML é uma forma de **meta-programação**

Programas que **criam programas** (hipóteses g) especializados em realizar a transformação *input-output*

- **Generalização**

O programa criado deve funcionar para exemplos previamente **não vistos**

A **teoria de ML** apresenta alguns resultados interessantes sobre isso

- **Complexidade da família de hipóteses** \times **quantidade de dados**

Quanto **mais complexos** forem as hipóteses, **mais dados** de treinamento são necessários

Supervisionado: Rótulos numéricos, y , representando:

- um valor numérico ($y \in \mathbb{R}$)
 - problema de regressão
 - estimar *redshift*
- um valor categórico ($y \in \{1, 2, \dots, c\}$)
 - problema de classificação
 - morfologia da galáxia ($y \in \{1, 2, \dots, c\}$)

Não-supervisionado: Y “desconhecido”

- clustering, etc

Exemplo

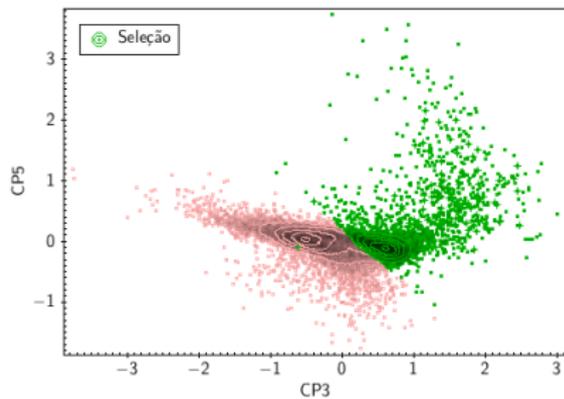
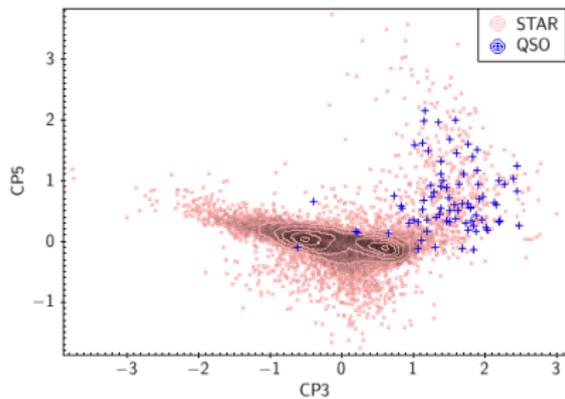
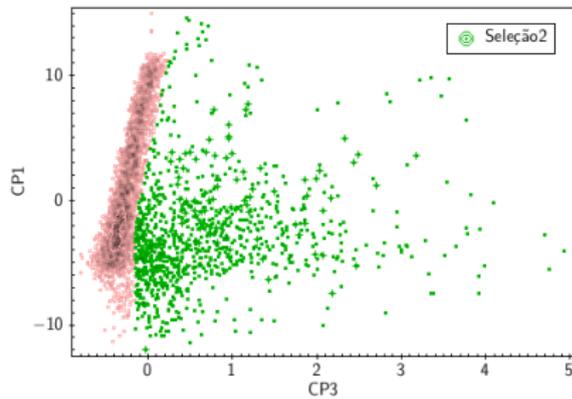
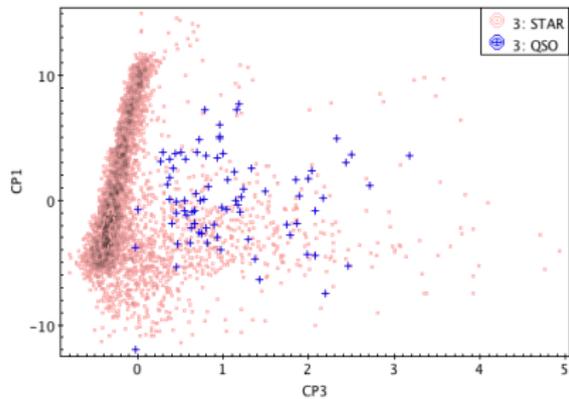
Deteção de quasares em catálogos

Detecção de quasares

(com Cláudia Oliveira, João Steiner, Liliane Nakazono)

- Dados: todos os objetos correspondentes a fontes pontuais no catálogo (ALHAMBRA)
- Como detectar fonte pontual que é QSO (e não estrela) ?
- Selecionar do catálogo os QSOs conhecidos (buscar QSO em listas de objetos identificados e fazer *match* de coordenada com os objetos do catálogo em estudo)
- Representação adequada para cada objeto do catálogo: *features* podem ser calculadas a partir da magnitude dos filtros
- Aplicar PCA (análise de componentes principais)
- Projecção 2D dos objetos – pares de componentes do PCA
- Selecionar projecção na qual QSOs (inicialmente selecionados) ficam concentrados em uma região compacta, mesmo que misturados às supostas estrelas
- Para cada objeto na região, calcular os k-NN em um conjunto conhecido de QSOs (possivelmente externos ao catálogo) – *redshift* homogêneo entre k-vizinhos confere maior confiança de que de fato aquele é um QSO, e de quebra fornece uma estimativa do seu *redshift*

Deteção de quasares



Detecção de quasares

(com Cláudia Oliveira, João Steiner, Liliane Nakazono)

- Dados: todos os objetos correspondentes a fontes pontuais no catálogo (ALHAMBRA)
- Como detectar fonte pontual que é QSO (e não estrela) ?
- Selecionar do catálogo os QSOs conhecidos (buscar QSO em listas de objetos identificados e fazer *match* de coordenada com os objetos do catálogo em estudo)
- Representação adequada para cada objeto do catálogo: *features* podem ser calculadas a partir da magnitude dos filtros
- Aplicar PCA (análise de componentes principais)
- Projecção 2D dos objetos - pares de componentes do PCA
- Selecionar projecção na qual QSOs (inicialmente selecionados) ficam concentrados em uma região compacta, mesmo que misturados às supostas estrelas
- Para cada objeto na região, calcular os k-NN em um conjunto conhecido de QSOs (possivelmente externos ao catálogo) - *redshift* homogêneo entre k-vizinhos confere maior confiança de que de fato aquele é um QSO, e de quebra fornece uma estimativa do seu *redshift*

zero dados
regras

```
if age > 40:
    if is_home_owner:
        print("give a credit")
    else:
        if income > 5000:
            print("give a credit")
        else:
            print("to refuse")
else:
    if education == "university":
        print("...")
    else:
        print("...")
```

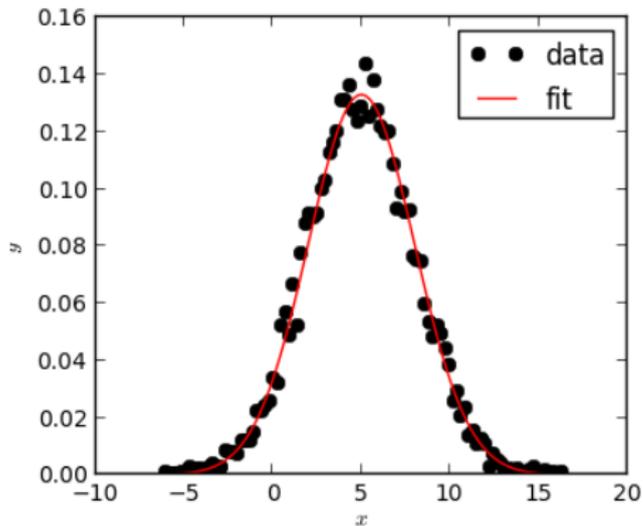
zero dados

regras

poucos dados

estimação paramétrica

template matching



Machine learning ao longo dos anos × dados

zero dados

regras

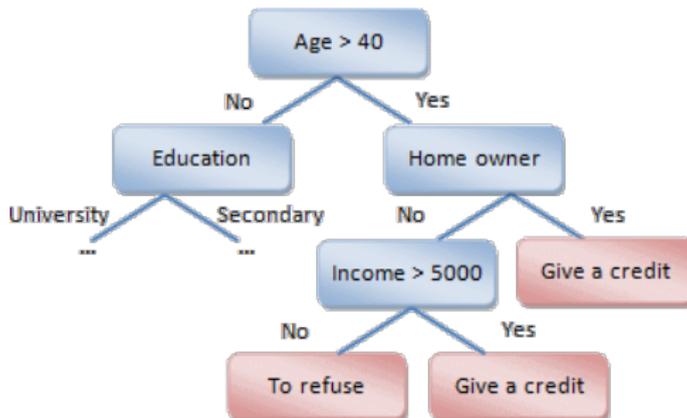
poucos dados

estimação paramétrica

template matching

mais dados

indução de modelos
(algoritmos de ML)



Machine learning ao longo dos anos × dados

zero dados

regras

poucos dados

estimação paramétrica

template matching

mais dados

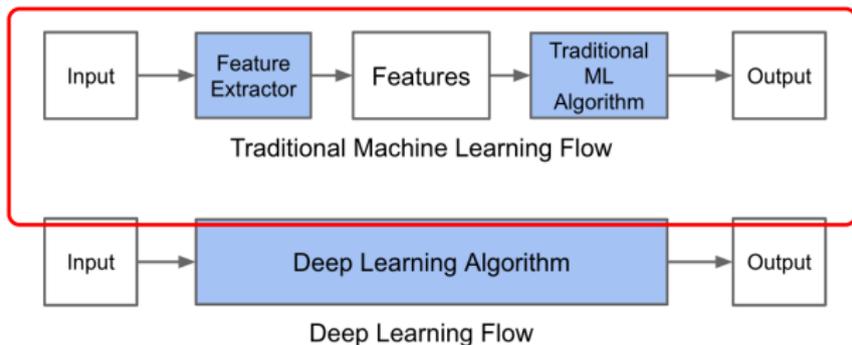
indução de modelos

(algoritmos de ML)

muito mais dados

engenharia de *features*

(representação de dados)



Machine learning ao longo dos anos × dados

zero dados

regras

poucos dados

estimação paramétrica

template matching

mais dados

indução de modelos

(algoritmos de ML)

muito mais dados

engenharia de *features*

(representação de dados)

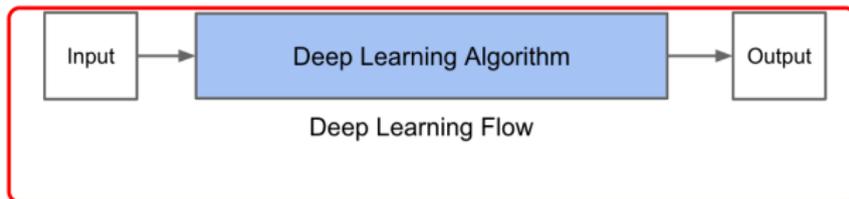
a era do Big Data

aprendizado de representação

(deep learning)



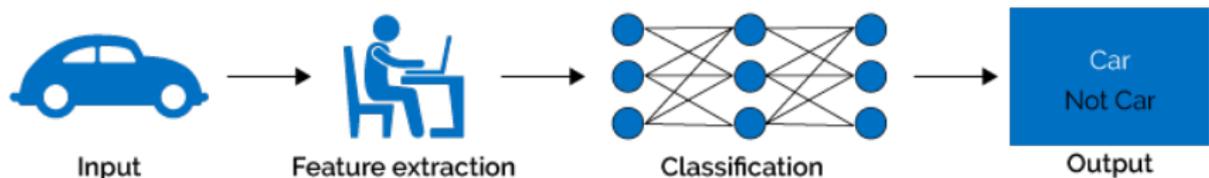
Traditional Machine Learning Flow



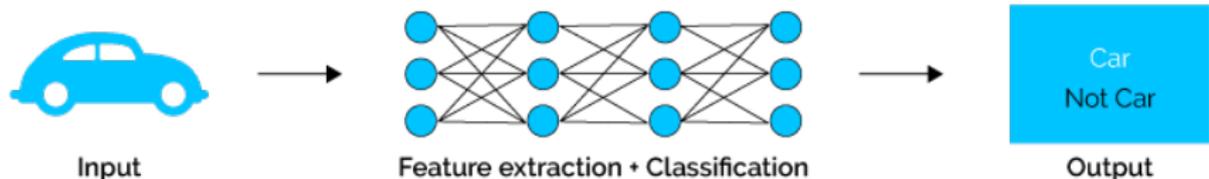
Deep Learning Flow

Traditional ML × Deep Learning

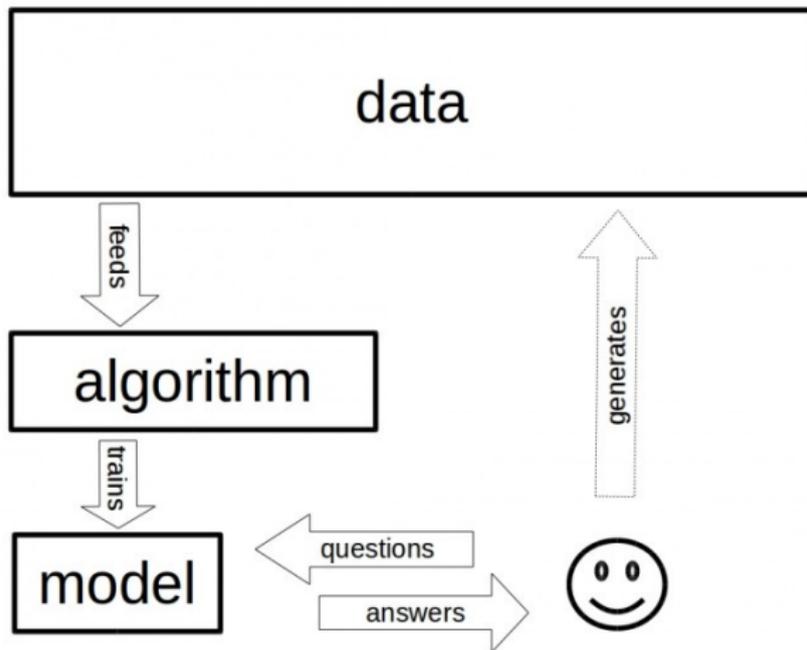
Machine Learning



Deep Learning



(<https://hackernoon.com/log-analytics-with-deep-learning-and-machine-learning-20a1891ff70e>)



<https://pieroit.org/portfolio/how-much-does-machine-learning-cost/>

Python, R : ambiente para prototipação

Scikit-learn : biblioteca de ferramentas úteis em ML

(<http://scikit-learn.org>)

TensorFlow, PyTorch : framework para desenvolvimento de *deep learning algorithms*

Intro a redes neurais: Livro do Nielsen

(<http://neuralnetworksanddeeplearning.com/>)

Intro a ML: curso da Caltech (<https://work.caltech.edu/telecourse.html>)

- Fizemos uma breve apresentação sobre **noções básicas de ML** (ponto de vista computacional)
- **Quando aplicar ML** : quando há alguma relação entrada-saída, difícil de ser descrita formalmente
- **Falta** : trabalhar o *gap* entre diferentes áreas
- **Utilidade bem pé no chão** : automatizar tarefas de “baixo-nível”, para dedicar mais tempo às tarefas de “alto-nível”