

# Análise de algoritmos de reconhecimento de padrões

Antônio Adelino da S. Neto<sup>a</sup>, Armstrong Lohãns de M. G. Quintino<sup>b</sup>

*Garanhuns, Brasil*

<sup>a</sup>*antonio.asn03@gmail.com*

<sup>b</sup>*lohansdemelo1108@gmail.com*

---

## Abstract

O objetivo principal desse trabalho foi a elaboração e o estudo de dois algoritmos para sistemas de aprendizado de máquina, os quais trabalhariam na classificação de textos em linguagem natural. Tais programas foram o algoritmo da Árvore de Decisão e o algoritmo de Naive Bayes. Inicialmente são mostrados os conceitos básicos sobre cada algoritmo de decisão acima citado. Depois haverá uma apresentação das técnicas usadas para a análise e as devidas conclusões.

*Keywords:* Reconhecimento de Padrões, Árvore de Decisão, Naive Bayes

---

## 1. Introdução

Os seres humanos e alguns outros animais possuem, entre outras habilidades, a aptidão no reconhecimento de padrões. O ser humano, especificamente, possui essa capacidade muito bem desenvolvida e tem uma enorme facilidade no reconhecimento formas, dando a elas significado e valor. Dentre elas pode-se citar a fisionomia de outros seres humanos, formas animais e vegetais, características pessoais e afins.

Essa habilidade sempre foi muito importante, pois foi por meio dela que a espécie humana conseguiu desenvolver-se com mais facilidade ao longo do tempo, uma vez que ela permite a assimilação e inferência de características em formas aparentemente semelhantes. Partindo dessa premissa, é possível notar a relevância dessa aptidão em reconhecimento para o ser humano, em especial o reconhecimento de padrões, visto que é por meio dela que consegue-se inferir em formas desconhecidas julgamentos prévios a partir de conhecimentos anteriores.

Com isso, afirma-se então que toda e qualquer forma de reconhecimento de padrões, por indivíduos, dá-se a partir de uma experiência passada. Dessa maneira é possível perceber que a destreza, ou não, no reconhecimento de padrões está

6 de outubro de 2019

18 diretamente vinculada aos estímulos que cada indivíduo foi submetido ao longo  
19 de sua vida [1].

20 Partindo dessas afirmativas, o presente artigo expõe um estudo que busca a  
21 análise comparativa de dois algoritmos. O algoritmo da Árvore de Decisão e o  
22 algoritmo de Naive Bayes, ambos voltados a classificação de dados baseando-se  
23 nos princípios de aprendizagem de máquina.

24 A análise desses algoritmos dá-se por meio da classificação de textos em lin-  
25 guagem natural, onde os classificadores recebem os textos e inferem a eles o sen-  
26 tido do que está escrito de acordo com classes anteriormente estabelecidas.

## 27 **2. Referencial Teórico**

### 28 *2.1. Algoritmo da Árvore de Decisão*

29 As Árvores de Decisão são técnicas muito populares de aprendizado de má-  
30 quina, são aplicadas às tarefas de classificação e regressão. Esta técnica é carac-  
31 terizada pelo seu modelo resultante, o qual é codificado como uma estrutura em  
32 árvore [2].

33 As árvores de decisão são algoritmos que buscam a classificação dos dados a  
34 partir da estruturação em árvore. O algoritmo divide um conjunto de dados em  
35 subconjuntos menores. Sabendo que o código estrutura-se em árvore, cada nó  
36 folha representa uma decisão.

37 Para chegar em uma decisão, o algoritmo comporta-se da seguinte maneira,  
38 com base nos valores dos recursos das instâncias, as árvores de decisão classificam  
39 os dados. Cada nó representa um recurso em uma instância da árvore de decisão  
40 que deve ser classificada, e cada ramo representa um valor [3].

41 Sabendo disso, percebe-se que cada dado, para ser classificado, passa por um  
42 conjunto finito de nós, tal conjunto é definido como as regras de classificação, pois  
43 a partir desse conjunto é possível saber o passo a passo do algoritmo, mostrando  
44 assim todas as regras que levaram a classificação daquela única instância. A prin-  
45 cipal vantagem do uso das árvores de decisão está justamente na capacidade do  
46 retorno dos passos para a decisão e não unicamente no resultado da classificação.

### 47 *2.2. Algoritmo de Naive Bayes*

48 Além das árvores de decisão, pode-se também fazer uso de outros tipos de  
49 classificadores, entre eles destaca-se o Naive Bayes, o qual possui uma análise  
50 dos dados a partir de conceitos probabilísticos, diferenciando-se das árvores de  
51 decisão.

52 O algoritmo de Naive Bayes é um classificador probabilísticos simples (ba-  
53 seado no Teorema de Bayes), tem como base em uma suposição comum de que  
54 todos os recursos são independentes um do outro [4]. A partir disso, ele descon-  
55 sidera completamente a correlação entre todas as variáveis, tratando cada variável  
56 de forma independente, esse algoritmo é frequentemente aplicado em processa-  
57 mento de linguagem natural.

58 Uma das principais vantagens do classificador Naive Bayes é que ele requer  
59 apenas uma pequena quantidade de dados iniciais de treinamento para poder esti-  
60 mar as médias e variações das variáveis necessárias para classificação [5].

### 61 **3. Ferramentas usadas**

#### 62 *3.1. Linguagem Python*

63 Antes de iniciar a implementação optamos por usar a linguagem de progra-  
64 mação Python. Essa linguagem, além de ser uma das mais populares do mundo,  
65 é vastamente usada na produção de softwares e algoritmos voltados aos concei-  
66 tos de aprendizagem de máquina, tanto no meio acadêmico quanto na indústria.  
67 Pode-se afirmar também que a linguagem de programação Python é uma lingua-  
68 gem de fácil manipulação e de grande eficiência produtiva o que facilita todo o  
69 processo de codificação.

#### 70 *3.2. Módulo Sickit-Learn*

71 Partindo dessa escolha inicial, foi possível utilizar o módulo Sickit-Learn para  
72 Python. Tal módulo integra em si uma grande quantidade de algoritmos de apren-  
73 dizagem de máquina que são voltados para problemas supervisionados e não su-  
74 pervisionados [6].

75 Além disso, essa biblioteca Python tem como característica a simplificação  
76 do uso de algoritmos de aprendizagem de máquina buscando a sua popularização.  
77 Essa difusão dá-se por meio da grande facilidade de uso, do bom desempenho e da  
78 sua documentação detalhada. Esse módulo ainda procura incentivar o seu uso, por  
79 meio de dependências mínimas e licença simplificada, em ambientes acadêmicos  
80 e comerciais de todo o mundo [6].

81 Diante dessas características encontradas na linguagem de programação Python  
82 e na biblioteca Sickit-Learn a implementação dos algoritmos, que servem de ex-  
83 perimento para o presente artigo, tornou-se mais rápida e mais objetiva.

### 84 3.3. *Wisdom Quotes (Base de dados)*

85 Para a criação da base de dados inicial dos classificadores foi necessário a  
86 utilização do site <http://wisdomquotes.com> a fim de escolher textos aleató-  
87 rios que fossem divididos por classificadores. O Wisdom Quotes conta com uma  
88 grade variedade de citações separadas por temas, permitindo a relação intuitiva  
89 citação-tema, o que vem a facilitar a criação e manipulação da base de dados dos  
90 classificadores.

## 91 4. Algoritmo e testes

### 92 4.1. *Criação da base de dados*

93 Antes de iniciar a implementação dos algoritmos brutos de aprendizado de  
94 máquina (Árvore de Decisão e Naive Bayers) tivemos que coletar as citações do  
95 site Wisdom Quotes. Para isso implementamos um script simples em Python  
96 (Crawler) que tinha como finalidade apenas a coleta das citações do endereço  
97 URL que passamos e colocar em arquivos de texto. Com o Crawler o processo de  
98 coleta de dados tornou-se muito mais rápido e eficaz. A coleta foi feita em três  
99 páginas do mesmo site, cada página continha citações de uma única categoria, as  
100 três categorias foram escolhidas randomicamente e são elas:

- 101 1. Peace (do inglês, "Paz")
- 102 2. Success (do inglês, "Sucesso")
- 103 3. Silence (do inglês, "Silêncio")

104 Com as citações salvas em arquivos de texto, tivemos que escolher aleatori-  
105 amente cento e cinquenta de cada classificador, a fim de deixar a base de dados  
106 uniforme, totalizando quatrocentos e cinquenta frases.

107 A partir disso, armazenamos essas frases em dois vetores, o vetor de treino  
108 e o vetor de teste. Essas listas continham respectivamente sessenta por cento e  
109 quarenta por cento das citações, escolhidas aleatoriamente entre si. Essas porcen-  
110 tagens foram estabelecidas visando um melhor desempenho dos classificadores,  
111 uma vez que a quantidade de dados para treino é um pouco maior quando compa-  
112 rada a quantidade de dados para testes.

### 113 4.2. *Árvore de Decisão*

114 Com a base de dados inicial pronta, implementamos a árvore de decisão. Para  
115 usar a função para criação de uma árvore de decisão no Sickit-Learn é necessá-  
116 rio inicialmente vetorizar as frases de teste, a função (também do Sickit-Learn)

```
x = CountVectorizer()
treinoFrase = x.fit_transform(treinoFrase).toarray()
```

Figura 1: Vetorizando frases de teste

117 *.fit\_transform()* do *CountVextorizer()* vetoriza as frases em forma de matriz, além  
118 de contar a repetição de cada palavra na frase, trecho ilustrado na Figura 1.

119 Com as frases devidamente vetorizadas, criamos a árvore de decisão a partir  
120 da função *tree.DecisionTreeClassifier()* e a treinamos com o comando *.fit(parm1,*  
121 *parm2)*, como mostrado na Figura 2, passando como parâmetro as frases vetori-  
122 zadas e a lista de classificações. Cada frase vetorizada está localizada na posição  
123 de índice correspondente ao seu classificador que está na lista de classificadores.

```
arvore = tree.DecisionTreeClassifier()
arvore.fit(treinoFrase, treinoClass)
```

Figura 2: Instanciando e treinando árvore de decisão

124 Após a criação da árvore, colocamos a base de testes para validar as classifi-  
125 cações, a função *.score(parm1, parm2)* retorna o percentual de acerto dos testes  
126 passados no segundo parâmetro, no primeiro parâmetro está o conjunto contendo  
127 todas as palavras conhecidas pela árvore (*Bag of Words*), veja na Figura 3.

```
arvore.score(bagOfWords, testeClass)
```

Figura 3: Função que retorna a porcentagem de acertos da estrutura

128 Com a árvore de decisão montada também é possível obter as regras de clas-  
129 sificação de uma determinada frase passada, para isso usamos a função, também  
130 do do Sickit-learn, *.decision\_path()*, conforme mostra a Figura 4.

131 Com isso, terminamos a implementação do algoritmo da árvore de decisão e  
132 suas funcionalidades principais, dando destaque a possibilidade de listar o cami-  
133 nho de decisão de uma citação.

#### 134 4.3. Naive Bayers

- 135 • Bullet point one
- 136 • Bullet point two

```
arvore.decision_path(frase)
```

Figura 4: Função que retorna a as regras de decisão da árvore

- 137 1. Numbered list item one  
138 2. Numbered list item two

#### 139 4.4. Subsection One

140 Quisque elit ipsum, porttitor et imperdiet in, facilisis ac diam. Nunc facilisis  
141 interdum felis eget tincidunt. In condimentum fermentum leo, non consequat leo  
142 imperdiet pharetra. Fusce ac massa ipsum, vel convallis diam. Quisque eget turpis  
143 felis. Curabitur posuere, risus eu placerat porttitor, magna metus mollis ipsum, eu  
144 volutpat nisl erat ac justo. Nullam semper, mi at iaculis viverra, nunc velit iaculis  
145 nunc, eu tempor ligula eros in nulla. Aenean dapibus eleifend convallis. Cras ut  
146 libero tellus. Integer mollis eros eget risus malesuada fringilla mattis leo facilisis.  
147 Etiam interdum turpis eget odio ultricies sed convallis magna accumsan. Morbi in  
148 leo a mauris sollicitudin molestie at non nisl.

Treatments	Response 1	Response 2
Treatment 1	0.0003262	0.562
Treatment 2	0.0015681	0.910
Treatment 3	0.0009271	0.296

Tabela 1: Table caption

#### 149 4.5. Subsection Two

150 Donec eget ligula venenatis est posuere eleifend in sit amet diam. Vestibulum  
151 sollicitudin mauris ac augue blandit ultricies. Nulla facilisi. Etiam ut turpis nunc.  
152 Praesent leo orci, tincidunt vitae feugiat eu, feugiat a massa. Duis mauris ipsum,  
153 tempor vel condimentum nec, suscipit non mi. Fusce quis urna dictum felis po-  
154 suere sagittis ac sit amet erat. In in ultrices lectus. Nulla vitae ipsum lectus, a  
155 gravida erat. Etiam quam nisl, blandit ut porta in, accumsan a nibh. Phasellus  
156 sodales euismod dolor sit amet elementum. Phasellus varius placerat erat, nec  
157 gravida libero pellentesque id. Fusce nisi ante, euismod nec cursus at, suscipit a  
158 enim. Nulla facilisi.

159 Integer risus dui, condimentum et gravida vitae, adipiscing et enim. Aliquam  
160 erat volutpat. Pellentesque diam sapien, egestas eget gravida ut, tempor eu nulla.

161 Vestibulum mollis pretium lacus eget venenatis. Fusce gravida nisl quis est mo-  
162 lestie eu luctus ipsum pretium. Maecenas non eros lorem, vel adipiscing odio.  
163 Etiam dolor risus, mattis in pellentesque id, pellentesque eu nibh. Mauris nec ante  
164 at orci ultricies placerat ac non massa. Aenean imperdiet, ante eu sollicitudin ves-  
165 tibulum, dolor felis dapibus arcu, sit amet fermentum urna nibh sit amet mauris.  
166 Suspendisse adipiscing mollis dolor quis lobortis.

$$e = mc^2 \quad (1)$$

## 167 5. The Second Section

168 Reference to Section 1. Etiam congue sollicitudin diam non porttitor. Etiam  
169 turpis nulla, auctor a pretium non, luctus quis ipsum. Fusce pretium gravida libero  
170 non accumsan. Donec eget augue ut nulla placerat hendrerit ac ut mi. Phasellus  
171 euismod ornare mollis. Proin tempus fringilla ultricies. Donec pretium feugiat  
172 libero quis convallis. Nam interdum ante sed magna congue eu semper tellus  
173 sagittis. Curabitur eu augue elit.

174 Aenean eleifend purus et massa consequat facilisis. Etiam volutpat placerat  
175 dignissim. Ut nec nibh nulla. Aliquam erat volutpat. Nam at massa velit, eu  
176 malesuada augue. Maecenas sit amet nunc mauris. Maecenas eu ligula quis turpis  
177 molestie elementum nec at est. Sed adipiscing neque ac sapien viverra sit amet  
178 vestibulum arcu rhoncus.

179 Vivamus pharetra nibh in orci euismod congue. Pellentesque habitant morbi  
180 tristique senectus et netus et malesuada fames ac turpis egestas. Quisque lacus  
181 diam, congue vel laoreet id, iaculis eu sapien. In id risus ac leo pellentesque  
182 pellentesque et in dui. Etiam tincidunt quam ut ante vestibulum ultricies. Nam at  
183 rutrum lectus. Aenean non justo tortor, nec mattis justo. Aliquam erat volutpat.  
184 Nullam ac viverra augue. In tempus venenatis nibh quis semper. Maecenas ac nisl  
185 eu ligula dictum lobortis. Sed lacus ante, tempor eu dictum eu, accumsan in velit.  
186 Integer accumsan convallis porttitor. Maecenas pretium tincidunt metus sit amet  
187 gravida. Maecenas pretium blandit felis, ac interdum ante semper sed.

188 In auctor ultrices elit, vel feugiat ligula aliquam sed. Curabitur aliquam elit  
189 sed dui rhoncus consectetur. Cras elit ipsum, lobortis a tempor at, viverra vitae  
190 mi. Cras sed urna sed eros bibendum faucibus. Morbi vel leo orci, vel faucibus  
191 orci. Vivamus urna nisl, sodales vitae posuere in, tempus vel tellus. Donec magna  
192 est, luctus non commodo sit amet, placerat et enim.

193 **Referências**

- 194 [1] P. Prado, A. Monteiro, Pattern recognition algorithms 5 (2008).
- 195 [2] G. Nuti, L. A. J. Rugama, A.-I. Cross, A bayesian decision tree algorithm,  
196 stat 1050 (2019) 11.
- 197 [3] R. Pandya, J. Pandya, C5. 0 algorithm to improved decision tree with fea-  
198 ture selection and reduced error pruning, International Journal of Computer  
199 Applications 117 (2015) 18–21.
- 200 [4] S. Xu, Bayesian naïve bayes classifiers to text classification, Journal of Infor-  
201 mation Science 44 (2018) 48–59.
- 202 [5] S. Vijayarani, S. Dhayanand, Liver disease prediction using svm and naïve  
203 bayes algorithms, International Journal of Science, Engineering and Techno-  
204 logy Research (IJSETR) 4 (2015) 816–820.
- 205 [6] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel,  
206 M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,  
207 D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine  
208 learning in Python, Journal of Machine Learning Research 12 (2011) 2825–  
209 2830.