

# Aplicação do Algoritmo k-NN

Pablo Luiz Leon

*Santo André – SP, UFABC*

---

## Abstract

This work aims to show the implementation of the kNN algorithm (k Nearest Neighbors), to predict a value, based on a set of values (Database).

The objective of the work is to implement the kNN algorithm in two ways, sequentially and in parallel. To show the benefits of parallel processing. artifacts.

*Keywords:* knn, spark, python, big data, parallel programming

---

## 1. Algoritmo K-NN

O algoritmo k-NN é um algoritmo de classificação o K vizinhos mais próximos (do inglês: K nearest neighbors – KNN). O KNN foi proposto por Fukunaga e Narendra em 1975 [1]. É um dos classificadores mais simples de ser implementado, de fácil compreensão e ainda hoje pode obter bons resultados dependendo de sua aplicação.

A ideia principal do KNN é determinar o rótulo de classificação de uma amostra baseado nas amostras vizinhas advindas de um conjunto de treinamento indicado.

Para tornar o entendimento mais fácil vamos aplicar um exemplo para explicar o seu funcionamento do algoritmo como o da Figura 1, na qual temos um problema de classificação com dois rótulos de classe e com  $k = 7$ . No exemplo, são aferidas as distâncias de uma nova amostra, representada por uma estrela, às demais amostras de treinamento, representadas pelas bolinhas azuis e amarelas. A variável  $k$  representa a quantidade de vizinhos mais próximos que serão utilizados para averiguar de qual classe a nova amostra pertence. Com isso, das sete amostras de treinamento mais próximas da nova amostra, 4 são do rótulo A e 3 do rótulo B. Portanto, como existem mais vizinhos do rótulo A, a nova amostra receberá o mesmo rótulo deles, ou seja, A.

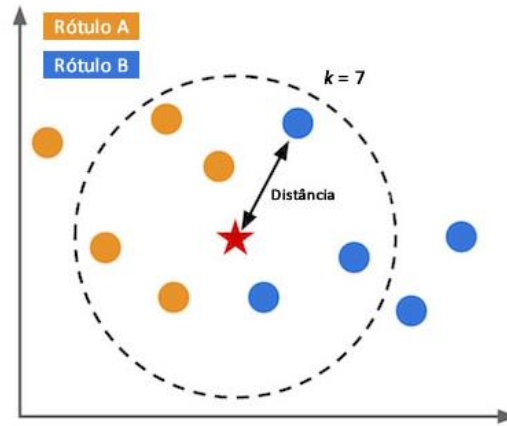


Figura 1: Explicação gráfica do Algoritmo k-NN

19 Dois pontos chave que devem ser determinados para aplicação do KNN são: a  
 20 métrica de distância e o valor de k. Para métrica de distância a mais utilizada é a  
 21 distância Euclidiana, descrita por:

$$D = \sqrt{(p_1 - q_1)^2 + \dots + (p_n - q_n)^2} = \sqrt{\sum_{i=1}^n (p_i - q_i)^2} \quad (1)$$

22 onde  $P = (p_1, \dots, p_n)$  e  $Q = (q_1, \dots, q_n)$  são dois pontos  $n$ -dimensionais. No  
 23 exemplo que citei da figura, essa distância seria calculada entre as bolinhas (azuis e  
 24 laranjas) e a estrela (a nova entrada). Como o exemplo usa dois parâmetros, cada  
 25 ponto tem dois valores  $x_1$  e  $x_2$ . Portanto para esse cenário do exemplo seria  $(x_1 - y_1)^2$   
 26  $+ (x_2 - y_2)^2$ , caso o seu problema tenha mais de dois parâmetros a abordagem é a  
 27 mesma para o cálculo euclidiano das distâncias.

28 Em relação ao valor k, não existe um valor único para a constante, a mesma  
 29 varia de acordo com a base de dados. É recomendável sempre utilizar valores  
 30 ímpares/primos, mas o valor ótimo varia de base para base, onde o sugerido para  
 31 esse algoritmo seria calcular o valor de k baseado no logaritmo na base dois do total  
 32 de amostras da base de dados em questão.

33 Em alto nível a grande vantagem do método do kNN é a sua abordagem simples  
 34 para ser compreendida e implementada, todavia, efetuar a distâncias dos pontos gerar  
 35 um custo computacional para calcular as distâncias de forma bem relevante no tempo

36 total da execução do algoritmo, caso tenha uma quantidade grande de amostras de  
37 treino essa atividade custosa pode gerar um tempo relativamente grande para a  
38 execução do algoritmo. Além disso, o método é sensível à escolha do k (quantidade  
39 de vizinhos observados).

---

```
1 inicialização:
2   Preparar conjunto de dados de entrada e saída
3   Informar o valor de k;
4 para cada nova amostra faça
5   Calcular distância para todas as amostras
6   Determinar o conjunto das k's distâncias mais próximas
7   O rótulo com mais representantes no conjunto dos k's
8   vizinhos será o escolhido
9 fim para
10 retornar: conjunto de rótulos de classificação
```

---

Figura 2: Passos do Algoritmo k-NN

## 40 2. Aplicação do para o trabalho proposto

41 Para um conjunto de dados onde temos os dados de cotação de ações da bolsa  
42 americana em um arquivo com o seguinte layout (x, y) onde:

43  $x \rightarrow$  data do valor da ação;

44  $y \rightarrow$  valor da ação.

45 Crie uma massa de dados onde dado um valor N determinado pelo usuário seja  
46 construído a massa de dados com o seguinte layout (l, v) onde:

47  $l \rightarrow$  lista de valores do ativo, onde total da lista é igual ao N especificado;

48  $v \rightarrow$  próximo valor da lista, na sequência de valores de ativo.

49 O Algoritmo foi dividido em algumas partes, funções que fazem parte específica  
50 do contexto todo do problema.

51 A Função criarDados(RDD\_Base, tamanho\_list), que efetua o preparo dos dados.  
52 A ideia é ter um arquivo que indica os valores de cada ação. Essa função recebe o  
53 arquivo e transforma os dados para que seja criado o RDD com os dados formatados  
54 para o cálculo do k-NN.

55 A Função DistanciaEuclidiana(RDDDados, DadosBase), que efetua o cálculo  
56 euclidiano dos valores indicado, onde recebe um RDD com os dados de treino, um  
57 valor de base para predição.

58 A Função `kNN(k, RDDTreino, Dados)`, que efetua o cálculo do k-NN, aonde  
59 passando a base de dados, o valor de k e a base de treino ele efetua o cálculo e retorna  
60 à predição do valor em questão para aquela sequência de dados.

61 Para esse trabalho temos a problemática de prever o valor do ativo dado a lista  
62 de dados com os valores do ativo indicado em uma lista. Portanto, para efetuar essa  
63 predição iremos pegar o k mais próximo do valor da lista de valores de ação em  
64 questão analisado, com isso o parâmetro k da implementação fica fixo em  $k=1$ .

65 Durante a execução do Algoritmo foram feitos testes para medir o tempo de  
66 execução, levando em consideração 1: modelo sequencial e 2: modelo paralelizado.

67 Para o modelo 1: sequencial, a carga no Data set de dados (RDD) inicialmente é  
68 feita de forma normal sem usar paralelismo, conforme o código abaixo:

```
137  
138 arquivo = os.path.join('Data', 'C:/UFABC/DADOS/ETFs/Dados_Testes.txt')  
139  
140 # Lê o arquivo com os dados e carrega em um RDD  
141 DadosRDD1 = (sc.textFile(arquivo, 8).collect())  
142
```

Figura 3: Código de Carregamento do RDD para execução em modo Sequencial

69 Com esse modelo adotado os dados de teste, apresentou um tempo total de 2  
70 minutos e 54 segundos para carregar os dados e efetuar o processamento das  
71 transformações.

72 Já no segundo modelo 2: modelo paralelizado, foi configurado na carga do Data  
73 set de dados (RDD) o início do paralelismo de processamento, conforme o código  
74 abaixo:

```
137  
138 arquivo = os.path.join('Data', 'C:/UFABC/DADOS/ETFs/Dados_Testes.txt')  
139  
140 # Lê o arquivo com os dados e carrega em um RDD  
141 DadosRDD1 = (sc.textFile(arquivo, 8).collect())  
142  
143 DadosRDD = sc.parallelize(DadosRDD1, 6)
```

Figura 4: Código de Carregamento do RDD para execução em modo Paralelo

75 Com esse modelo os testes apresentaram uma melhora significativa no tempo  
76 aonde levou cerca de 1 Minuto e 52 segundos para carregar os dados e efetuar o  
77 processamento das transformações. Gerando no total uma eficiência de 36% do  
78 tempo de processamento.

## **Referências**

[1] Fukunaga, K.; Narendra, P. M. A branch and bound algorithm for computing k-nearest neighbors. *IEEE Transactions on Computers*, v. 100, n. 7, p. 750–753, 1975.