

Algorytm k-NN

Filip Galas
Dominik Góralczyk

Agenda

1. Wstęp teoretyczny
 - a. intuicja, algorytm k-Nearest Neighbors
 - b. krótki przegląd kilku algorytmów Nearest Neighbor Search
 - c. kernel k-NN
 - d. metryki
2. Przykłady w Pythonie (<https://github.com/zalon525/mro-knn>)

Algorytm

Mamy zbiór treningowy złożony z obiektów (wektorów), dla których znamy ich klasy (w przypadku klasyfikacji) lub wartości (w przypadku regresji).

1. Dla danego punktu testowego znajdujemy k najbliższych sąsiadów w zbiorze treningowym.
2. W zależności od rodzaju problemu:
 - a. W przypadku klasyfikacji przydzielamy punkt testowy do klasy, do której należy większość k najbliższych sąsiadów (tzw. głosowanie)
 - b. W przypadku regresji wartość punktu testowego jest średnią wartości k najbliższych sąsiadów.

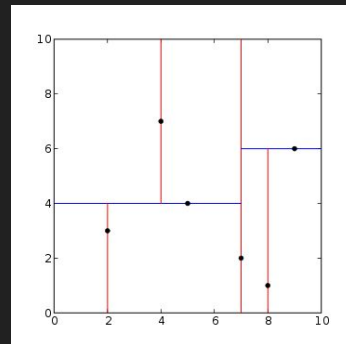
k - jest parametrem algorytmu (jest to najczęściej niewielka liczba naturalna)

Jak znajdować najbliższych sąsiadów?

Jest kilka możliwości:

1. Metody dokładne:

- Podjęście naiwne - dla danego punktu przeszukujemy cały zbiór, aby znaleźć najbliższego sąsiada - złożoność $O(n)$
- Podział przestrzeni - metody polegające na podziale przestrzeni, tak aby ograniczyć przeszukiwany obszar. Np. istnieje metoda wykorzystująca drzewo k-d (zbudowanie drzewa k-d - $O(n \log n)$, znalezienie najbliższego sąsiada - średnio $O(\log n)$)



Jak znajdować najbliższych sąsiadów

2. Metody aproksymacyjne (parę przykładów):

- Locality-sensitive hashing (LSH) - metoda polegająca na przydzielaniu punktów do “kubeków” wg funkcji hashującej skonstruowanej tak, aby dla punktów sobie bliskich zwracała tę samą wartość z dużym prawdopodobieństwem
- Cover tree

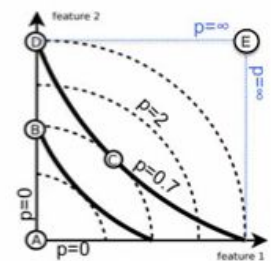
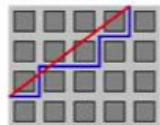
kernel k-NN

$$\max_r \left(\sum_{i=1}^k K(D(x, x_{(i)})) I(y_{(i)} = r) \right) .$$

- rectangular kernel $\frac{1}{2} \cdot I(|d| \leq 1)$
- triangular kernel $(1 - |d|) \cdot I(|d| \leq 1)$
- Epanechnikov kernel $\frac{3}{4}(1 - d^2) \cdot I(|d| \leq 1)$
- quartic or biweight kernel $\frac{15}{16}(1 - d^2)^2 \cdot I(|d| \leq 1)$

• Minkowski distance (p -norm): $D(x, x') = \sqrt[p]{\sum_d |x_d - x'_d|^p}$

- $p=2$: Euclidian
- $p=1$: Manhattan
- $p=0$: Hamming ... logical AND
- $p=\infty$: $\max_d |x_d - x'_d|$... logical OR



Metrics intended for real-valued vector spaces:

identifier	class name	args	distance function
"euclidean"	EuclideanDistance	•	<code>sqrt(sum((x - y)^2))</code>
"manhattan"	ManhattanDistance	•	<code>sum(x - y)</code>
"chebyshev"	ChebyshevDistance	•	<code>max(x - y)</code>
"minkowski"	MinkowskiDistance	p	<code>sum(x - y ^p)^(1/p)</code>
"wminkowski"	WMinkowskiDistance	p, w	<code>sum(w * x - y ^p)^(1/p)</code>
"seuclidean"	SEuclideanDistance	V	<code>sqrt(sum((x - y)^2 / V))</code>
"mahalanobis"	MahalanobisDistance	V or VI	<code>sqrt((x - y)' V^-1 (x - y))</code>

Metrics intended for two-dimensional vector spaces: Note that the haversine distance metric requires data in the form of [latitude, longitude] and both inputs and outputs are in units of radians.

identifier	class name	distance function
"haversine"	HaversineDistance	<code>2 arcsin(sqrt(sin^2(0.5*dx)</code> <div>• <code>cos(x1)cos(x2)sin^2(0.5*dy))</code></div>

Metrics intended for integer-valued vector spaces: Though intended for integer-valued vectors, these are also valid metrics in the case of real-valued vectors.

identifier	class name	distance function
"hamming"	HammingDistance	<code>N_unequal(x, y) / N_tot</code>
"canberra"	CanberraDistance	<code>sum(x - y / (x + y))</code>
"braycurtis"	BrayCurtisDistance	<code>sum(x - y) / (sum(x) + sum(y))</code>