

UNIVERSIDAD NACIONAL DE AGUSTÍN
FACULTAD DE INGENIERÍA DE PRODUCCIÓN Y
SERVICIOS

ESCUELA PROFESIONAL DE CIENCIA DE LA COMPUTACIÓN



Kd Tree - KNN

Alumno:

HERMOZA LOAYZA, MIGUEL
ÁNGEL
HERRERA COOPER, MIGUEL
ALEXANDER

Docente:

Machaca Arceda
VICENTE

Arequipa - Perú

Índice

1. Resumen	2
2. Introducción	2
3. Range Query	3
3.1. Ejercicio 1	4
3.2. Ejercicio 2	5
3.3. Resultados	6
4. Knn	7
4.1. Definición	7
4.2. ¿Cómo funciona?	7
5. Imagen intuición de clasificación con KNN	10
5.1. ¿Cómo puede una función asociar una etiqueta con una forma?	12
5.2. Preparando los datos	12
5.3. Algoritmo	13
5.4. Conjunto de entrenamiento, conjunto de validación, conjunto de prueba	16
6. Aplicaciones	16
7. Conclusión	17
8. Bibliografía	17

Índice de figuras

1. Resumen

En este documento presentamos la implementación función Range Query y haremos uso del Algoritmo Knn para clasificar imágenes.

Presentamos un algoritmo K-D Tree es una estructura de datos espacial común en muchos métodos de análisis de datos, se pueden ejecutar cuando se ha indexado todo el conjunto de datos, es decir, no están en línea. Incluso las pocas implementaciones en línea no son progresivas en el sentido de que el tiempo para indexar los datos entrantes no está limitado y puede exceder la latencia requerida por los sistemas progresivos. Esta latencia restringe significativamente la interactividad de los sistemas de visualización, especialmente cuando se trata de datos a gran escala. Mejoramos los k-dtrees tradicionales para el vecino progresivo más cercano al más cercano, lo que permite consultas KNN rápidas e indexa continuamente nuevos lotes de datos cuando es necesario. Siguiendo el paradigma de computación progresiva, nuestro k-dtree progresivo está limitado en el tiempo, lo que permite a los analistas acceder a resultados continuos dentro de una latencia interactiva. También presentamos puntos de referencia de rendimiento para comparar en línea y progresivo-dtrees.

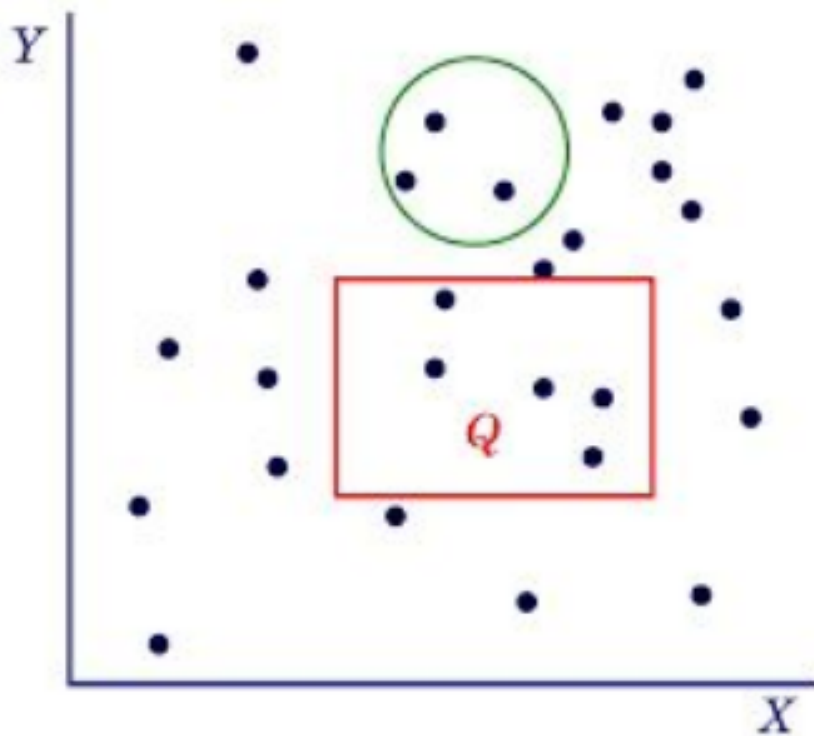
2. Introducción

La idea es realmente sencilla: el algoritmo clasifica cada dato nuevo en el grupo que corresponda, según tenga k vecinos más cerca de un grupo o de otro. Es decir, calcula la distancia del elemento nuevo a cada uno de los existentes, y ordena dichas distancias de menor a mayor para ir seleccionando el grupo al que pertenecer. Este grupo será, por tanto, el de mayor frecuencia con menores distancias.

El K-NN es un algoritmo de aprendizaje supervisado, es decir, que a partir de un juego de datos inicial su objetivo será el de clasificar correctamente todas las instancias nuevas. El juego de datos típico de este tipo de algoritmos está formado por varios atributos descriptivos y un solo atributo objetivo (también llamado clase).

3. Range Query

- La función RangeQuery , devuelve un conjunto de puntos a partir de una área de búsqueda, por ejemplo en la Figura de abajo, el círculo y el rectángulo representan el rango de búsqueda. La estructura debe ser capaz de retornar el conjunto de puntos dentro de dicho rango sin consultar todos los puntos.



3.1. Ejercicio 1

Implemente la función rangeQueryCircle en su kdtree:

```
1 function range_query_circle ( node , center , radio , queue ,
2     depth = 0 ) {
3     if ( node === null ) return null ;
4
5     var axis = node . axis ;
6     var nb = null ;
7     var ob = null ;
8
9     if ( center [ axis ] < node . point [ axis ] ) {
10         nb = node . left ;
11         ob = node . right ;
12     } else {
13         nb = node . right ;
14         ob = node . left ;
15     }
16
17     var best = closer_point ( center , node , range_query_circle ( nb , center
18         , radio , queue , depth + 1 ) ) ;
19
20     if ( Math . abs ( center [ axis ] - node . point [ axis ] ) <= radio ||
21         distanceSquared ( center , best . point ) > Math . abs ( center [ axis ] -
22         node . point [ axis ] ) ) {
23
24         if ( distanceSquared ( center , node . point ) <= radio ) {
25             queue . push ( node . point ) ;
26         }
27
28         best = closer_point ( center , best , range_query_circle ( ob , center ,
29             radio , queue , depth + 1 ) ) ;
30     }
31
32     return best ;
33 }
```

circleQ.js

3.2. Ejercicio 2

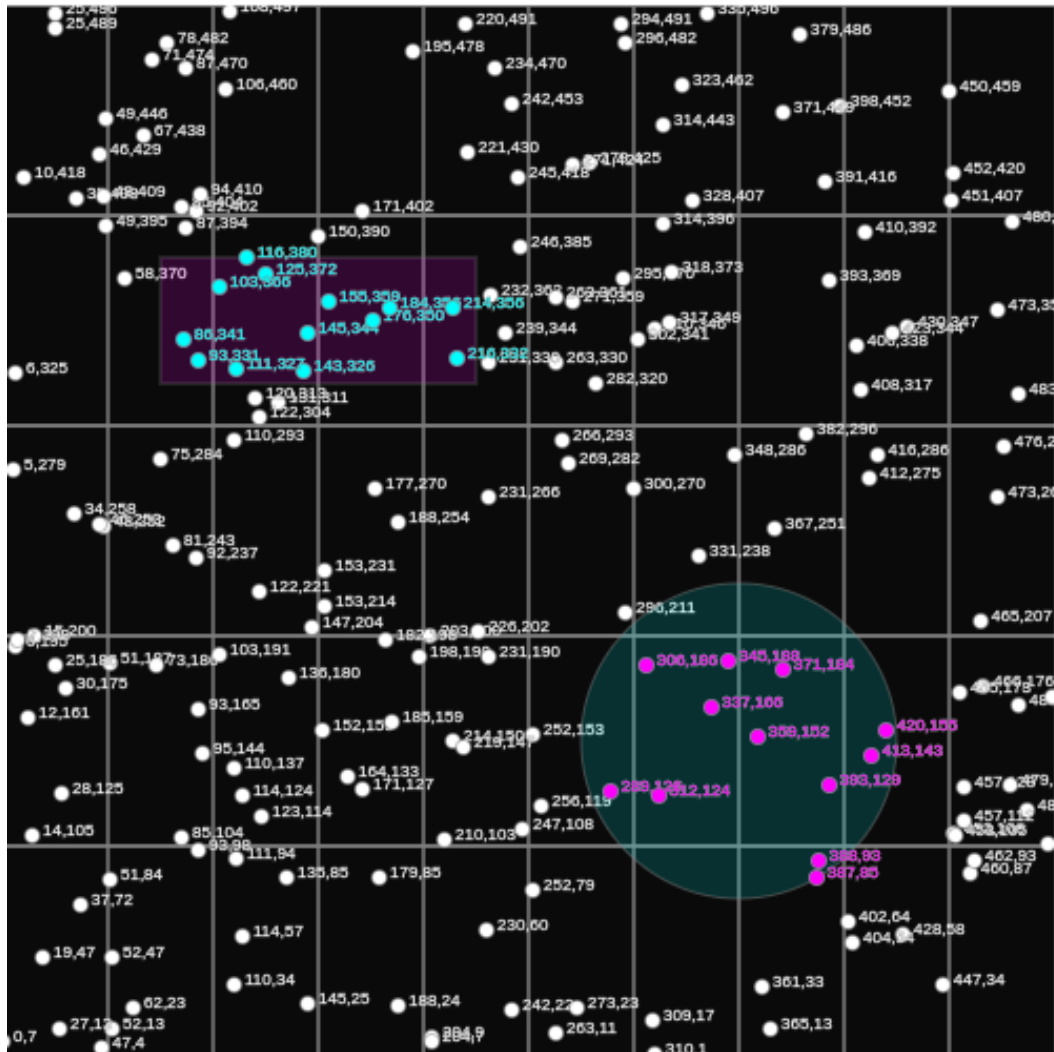
Implemente la función rangeQueryRectangle en su kdtree:

```
1 function range_query_rect ( node , center , hug , queue , depth
   = 0 ){
2   if (node==null) return null;
3
4   var axis = node.axis ;
5   var nb = null;
6   var ob = null;
7
8   if ( center [ axis]<node.point [ axis ]){
9     nb=node.left ;
10    ob=node.right ;
11  } else {
12    nb=node.right ;
13    ob=node.left ;
14  }
15  var best=closer_point ( center , node , range_query_rect (nb, center ,
    hug , queue , depth+1));
16
17  if (Math.abs ( center [ axis]-node.point [ axis ])<=hug [ axis]*2 ||
    distanceSquared ( center , best . point )>Math.abs ( center [ axis]-node
    . point [ axis ])) {
18
19    if (Math.abs ( center [0]-node.point [0])<=hug [0] && Math.abs (
    center [1]-node.point [1])<=hug [1]) {
20
21      queue.push ( node . point );
22    }
23    best=closer_point ( center , best , range_query_rect (ob, center , hug
    , queue , depth+1));
24  }
25
26  return best ;
27 }
```

rectangleQ.js

3.3. Resultados

- Usando la función rangeQueryCircle y rangeQueryRectangle obtene-
mos los siguientes resultados determinando el conjunto de puntos den-
tro del rango de búsqueda de cada estructura.



4. Knn

4.1. Definición

La idea básica sobre la que se fundamenta este paradigma es que un nuevo caso se va a clasificar en la clase mas frecuente a la que pertenecen sus K vecinos mas cercanos. El paradigma se fundamenta por tanto en una idea muy simple e intuitiva, lo que unido a su fácil implementación hace que sea un paradigma clasificadorio muy extendido.

Después de introducir el algoritmo K-NN básico y presentar algunas variantes del mismo, en este tema se estudian métodos para la selección de prototipos

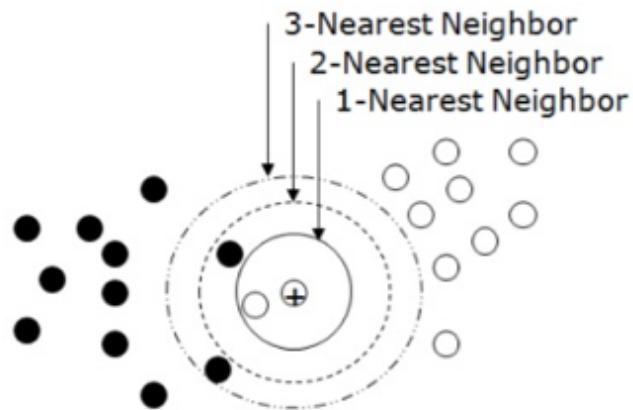
4.2. ¿Cómo funciona?

En contraste con otros algoritmos de aprendizaje supervisado, K-NN no genera un modelo fruto del aprendizaje con datos de entrenamiento, sino que el aprendizaje sucede en el mismo momento en el que se prueban los datos de test. A este tipo de algoritmos se les llama lazy learning methods.

Destacar que K-NN es muy sensible a:

- La variable k, de modo que con valores distintos de k podemos obtener resultados también muy distintos. Este valor suele fijarse tras un proceso de pruebas con varias instancias.
- La métrica de similitud utilizada, puesto que esta influirá, fuertemente, en las relaciones de cercanía que se irán estableciendo en el proceso de construcción del algoritmo. La métrica de distancia puede llegar a contener pesos que nos ayudarán a calibrar el algoritmo de clasificación, convirtiéndola, de hecho, en una métrica personalizada.

Vemos cómo influye el valor de k :



- Para $k = 1$ el algoritmo clasificará la bola con signo $+$ como blanca
- Para $k = 2$ el algoritmo no tiene criterio para clasificar la bola con signo $+$
- Para $k = 3$ el algoritmo clasificará la bola con signo $+$ como negra

Su mayor debilidad es la lentitud en el proceso de clasificación puesto que su objetivo no es obtener un modelo optimizado, sino que cada instancia de prueba es comparada contra todo el juego de datos de entrenamiento y, será la bondad de los resultados lo que determinará el ajuste de aspectos del algoritmo como el propio valor k , el criterio de selección de instancias para formar parte del juego de datos “D” de entrenamiento o la propia métrica de medida de similitud.

Finalmente, conoceremos la precisión del algoritmo a través de una evaluación de la matriz de confusión. En ella, se valorarán los verdaderos positivos, los falsos positivos, los positivos incorrectos, etc.

		Clase verdadera		
		C1	C2	
Clase predicha	C1	a	b	p_1
	C2	c	d	p_2
		π_1	π_2	

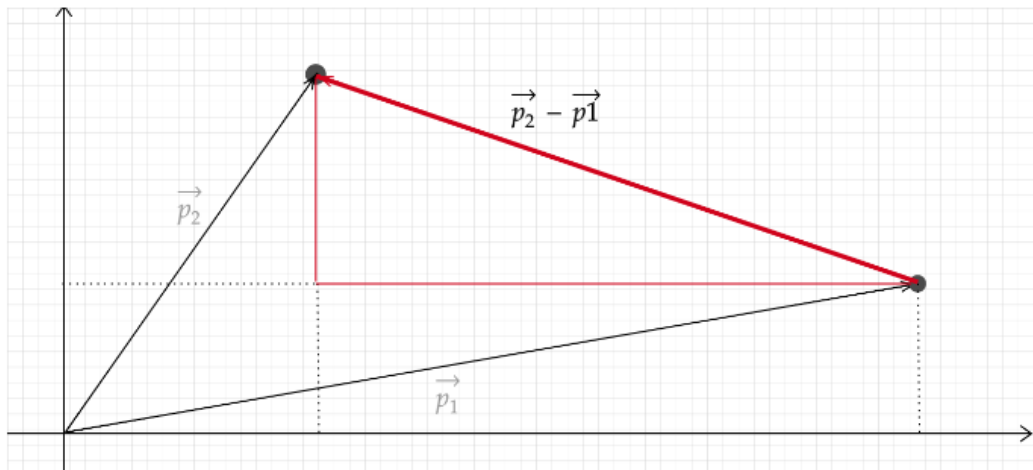
5. Imagen intuición de clasificación con KNN

Cada punto en el ejemplo de espacio 2D de KNN se puede representar como un vector (por ahora, una lista de dos números). Todos esos vectores apilados verticalmente formarán una matriz que representa todos los puntos en el plano 2D.

$$\begin{array}{cccc}
 P & C_x & C_y & color \\
 \vec{p}_1 & [10 & 20] & green \\
 \vec{p}_2 & [11 & 30] & orange \\
 \vdots & \vdots & \vdots & \vdots \\
 \vec{p}_n & [90 & 45] & green
 \end{array}
 \succ
 \begin{array}{cc}
 X & Y \\
 \begin{bmatrix} 10 & 20 \\ 11 & 30 \\ \vdots & \vdots \\ 90 & 40 \end{bmatrix} & \begin{bmatrix} green \\ orange \\ \vdots \\ green \end{bmatrix}
 \end{array}
 \succ
 \begin{array}{cc}
 X & Y \\
 \begin{bmatrix} \vec{p}_1 \\ \vec{p}_2 \\ \vdots \\ \vec{p}_n \end{bmatrix} & \begin{bmatrix} green \\ orange \\ \vdots \\ green \end{bmatrix}
 \end{array}$$

$C_x = x \text{ coordinate of a point}$
 $C_y = y \text{ coordinate of a point}$

En un plano 2D, si cada punto es un vector, la distancia euclidiana (escalar) se puede derivar del siguiente diagrama:



Esto, esencialmente, es una representación geométrica de la ecuación de la Distancia Euclidiana:

$$d(p_1, p_2) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \text{ where } p_1 = (x_1, y_1) \text{ and } p_2 = (x_2, y_2)$$

$$d(p_1, p_2) = \sqrt{(p_1^1 - p_2^1)^2 + (p_1^2 - p_2^2)^2}$$

Ahora piense en una imagen de 32 x 32 . Tenemos un total de 32 x 32 = 1024 píxeles.

Entonces, si tenemos que representar una imagen usando una estructura que la computadora pueda entender, tendremos un vector grande de 1024 elementos de largo.

Pero, ¿cuál será el valor de los números?

Por supuesto, el color que tiene el píxel.

El color de un píxel se puede representar utilizando 3 valores de color (supongamos que el canal de color RGB) por lo que un píxel contiene una lista de 3 valores (1 valor para r, 1 valor para g y 1 valor para b del canal RGB). Codificamos esto en una estructura similar, pero creamos 1024 valores de r, seguidos de 1024 valores de g, seguidos de 1024 valores de b. Eso es un total de 3072 valores en un vector.

Si no tenemos más de dos objetos, entonces tenemos una longitud de vector de 3072 apilados verticalmente. La matriz se vería así:

$$\begin{bmatrix} px_r^{1,1} & px_r^{1,2} & \dots & px_r^{32,32} & px_g^{1,1} & px_g^{1,2} & \dots & px_g^{32,32} & px_b^{1,1} & px_b^{1,2} & \dots & px_b^{32,32} \\ px_r^{1,1} & px_r^{1,2} & \dots & px_r^{32,32} & px_g^{1,1} & px_g^{1,2} & \dots & px_g^{32,32} & px_b^{1,1} & px_b^{1,2} & \dots & px_b^{32,32} \end{bmatrix}$$

5.1. ¿Cómo puede una función asociar una etiqueta con una forma?

Pasamos una forma codificada (gato, perro, etc.) y le decimos al programa que esto es 'gato', esto es 'perro'. El programa representa imágenes como puntos en el espacio dimensional superior y les asigna una etiqueta. Esto se llama entrenar al modelo .

Una vez que se completa el entrenamiento, pasamos los datos de prueba al modelo para predecir el animal, al igual que le dimos un nuevo punto en el espacio 2D para determinar su color.

5.2. Preparando los datos

La depuración de un modelo con entrenamiento real y datos de prueba suele ser tedioso. Por supuesto, crear mis propios datos no hará que este programa esté libre de errores, pero este enfoque garantiza que los errores rudimentarios se puedan encontrar de manera rápida y oportuna.

La idea es simple: crear centros de clúster en función del número de clases que desea crear. Genere aleatoriamente más puntos alrededor del centro del clúster, escalados por una desviación estándar predefinida.

El valor de retorno de esta función es (X, y). X contiene vectores apilados verticalmente e y es el color del punto. Cada artículo en la pila es un punto en el plano XY.

$$\begin{array}{cccc}
 \vec{P} & C_x & C_y & color \\
 \vec{p}_1 & [10 & 20] & green \\
 \vec{p}_2 & [11 & 30] & orange \\
 \vdots & \vdots & \vdots & \vdots \\
 \vec{p}_n & [90 & 45] & green
 \end{array}
 \succ
 \begin{array}{cc}
 X & Y \\
 \begin{bmatrix} 10 & 20 \\ 11 & 30 \\ \vdots & \vdots \\ 90 & 40 \end{bmatrix} & \begin{bmatrix} green \\ orange \\ \vdots \\ green \end{bmatrix}
 \end{array}
 \succ
 \begin{array}{cc}
 X & Y \\
 \begin{bmatrix} \vec{p}_1 \\ \vec{p}_2 \\ \vdots \\ \vec{p}_n \end{bmatrix} & \begin{bmatrix} green \\ orange \\ \vdots \\ green \end{bmatrix}
 \end{array}$$

$C_x = x \text{ coordinate of a point}$
 $C_y = y \text{ coordinate of a point}$

5.3. Algoritmo

Usamos matrices numpy para realizar el cálculo de distancia. Cuando un programa no es entendible, debemos ejecutar cada línea en Python REPL para ver cuál es la entrada y cuál es la salida. Esto ayuda a entender por qué se hizo este método en particular.

```
1 import networkx as nx
2 import matplotlib.pyplot as plt
3 import cv2
4 import math
5
6
7
8 k=3072
9
10 class Node():
11
12     def __init__(self, point):
13         self.point = point
14         self.axis = None
15         self.left = None
16         self.right = None
17         self.tipo = None
18
19 class Cola_prioridad():
20
21     def __init__(self, cant):
22         self.cantidad = cant
23         self.vector=[]
24
25     def insertar(self, distancia, x):
26         for i in range(len(self.vector)):
27             if (self.vector[i][1]==x):
28                 return
29
30         if (len(self.vector)==self.cantidad):
31             if (self.vector[self.cantidad-1][0]>distancia):
32                 self.vector[self.cantidad-1]=[distancia, x]
33                 self.vector.sort(key=lambda tup: tup[0])
34         else:
35             self.vector.append([distancia, x])
36             self.vector.sort(key=lambda tup: tup[0])
37
38     def top(self):
39         return self.vector[len(self.vector)-1][0]
40
41     def capacidad(self):
42         if (len(self.vector)==self.cantidad):
```

```

43         return True
44     else :
45         return False
46
47
48 def build_kdtree(points , depth=0):
49     if not points:
50         return None
51
52     axis = depth % k
53
54     points.sort(key=lambda tup: tup.point[ axis ])
55
56     median = len(points)//2
57
58     node = points[median]
59     node.axis=axis
60
61     node.left=build_kdtree(points[:median] , depth+1)
62     node.right=build_kdtree(points[median+1:] , depth+1)
63
64     return node
65
66 def distanceSquared(a, b):
67     distance = 0
68     for i in range(k):
69         distance = distance + pow((a[i]-b[i]),2)
70     return math.sqrt(distance)
71
72 def closest_point(node, point , depth , cola):
73     if (node == None):
74         return
75     axis = depth % k
76     nb = None
77     ob = None
78     cola.insertar(distanceSquared(point , node.point) , node)
79     if (point[ axis ] < node.point[ axis ]):
80         nb = node.left
81         ob = node.right
82     else:
83         nb = node.right
84         ob = node.left
85
86     closest_point(nb, point , depth+1, cola)
87
88     if( ( cola.capacidad()==False) or ( cola.top() > abs(point[ axis
89 ]-node.point[ axis ]))):
90         closest_point(ob, point , depth+1, cola)

```

```

91 def nearest(root, point, count):
92     ColaP = Cola_prioridad(count)
93     closest_point(root, point, 0, ColaP)
94     return ColaP.vector
95
96
97 def image_to_feature_vector(image, size=(32, 32)):
98     return cv2.resize(image, size).flatten()
99
100
101 def main():
102
103     puntos=[]
104     normal=0
105     trampa=0
106
107     for i in range(1,11):
108         direccion1="/home/maik/Escritorio/eda/kdtree/img/normal"
109         +str(i)+".png"
110         direccion1="/home/maik/Escritorio/eda/kdtree/img/trampa"
111         +str(i)+".png"
112         image1 = cv2.imread(direccion1)
113         image1 = cv2.imread(direccion1)
114         pixels1 = image_to_feature_vector(image1)
115         pixels1 = image_to_feature_vector(image1)
116         nodo1 = Node(pixels1)
117         nodo1 = Node(pixels1)
118         nodo1.tipo="normal"
119         nodo1.tipo="trampa"
120         puntos.append(nodo1)
121         puntos.append(nodo1)
122
123
124 root=build_kdtree(puntos)
125 image = cv2.imread("/home/maik/Escritorio/eda/kdtree/img/
126 test.png")
127 pixels = image_to_feature_vector(image)
128 resul = nearest(root, pixels,7)
129
130 for i in resul:
131     if(i[1].tipo=="normal"):
132         normal=normal+1
133     else:
134         trampa=trampa+1
135
136 if(trampa>normal):
137     print("Es carta Trampa")
138 else:

```



```
137     print("Es carta Monstruo")
138
139
140 main()
```

imagenes.py

5.4. Conjunto de entrenamiento, conjunto de validación, conjunto de prueba

Usamos una pequeña base de datos de imágenes, que es un pequeño conjunto de datos de 10 clases. Este conjunto de datos ya está segregado en conjunto de entrenamiento y conjunto de prueba. Sin embargo, deberíamos crear un conjunto de validación a partir del conjunto de entrenamiento usando la relación de corte 1: 9 (si el conjunto de entrenamiento tenía 1k filas después del corte, 100 filas van en el conjunto de validación 900 filas van en el conjunto de entrenamiento).

Y necesitamos este conjunto de validación para ajustar nuestro hiperparámetro (parámetro que no se deriva del proceso de aprendizaje, pero es parte del proceso de aprendizaje) - K.

Tenga en cuenta que he usado el conjunto de entrenamiento y el conjunto de prueba solo con fines ilustrativos en la base de código.

6. Aplicaciones

- Este modelo se utiliza en data mining (minería de datos), por lo que su uso en aprendizaje automático (Machine Learning) es una de sus aplicaciones clave para, por ejemplo, sistemas de recomendación (plataformas de contenido digital, procesos de venta cruzada en eCommerce, etc.).
- El concepto de los vecinos más cercanos es un importante factor en la física de estado sólido, ya que las propiedades de la materia se modifican en función de la estructura de las redes y las distancias entre los átomos. De hecho, es habitual hacer aproximaciones en los modelos físicos para despreciar aquellas partículas que no forman parte de los K vecinos más próximos.

7. Conclusión

La precisión del algoritmo no es muy eficiente. Lo interesante es que incluso la implementación de vanguardia del algoritmo produce resultados similares; claramente, todavía tenemos un largo camino por recorrer, una vez que nuestras bases están cubiertas.

8. Bibliografía

- <https://www.analiticaweb.es/algoritmo-knn-modelado-datos/>
- https://es.wikipedia.org/wiki/K_vecinos_m
- <https://medium.com/@YearsOfNoLight/intro-to-image-classification-with-knn-987bc112f0c2>
- <https://www.pyimagesearch.com/2016/08/08/k-nn-classifier-for-image-classification/>