

# Descripción del Algoritmo

En este documento se presenta la información necesaria para implementar el algoritmo de búsqueda de aplicaciones móviles más parecidas a otra, para esto se usó la estructura de datos espacial **k-dimensional tree**

## K-dimensional Tree

Esta implementación se hizo en 2 pasos:

1. **Inserción:** Se leyó las filas del csv, y se guardaron en un dataframe que facilita la librería pandas de python. Luego se crea el árbol ordenándolo por las dimensiones identificadas y cada fila del dataframe se guardó en la “clase nodo”
  - 1.1 Los datos identificados como dimensiones son los siguientes:
    - **Id**
    - **size\_bytes**
    - **Price**
    - **rating\_count\_tot**
    - **rating\_count\_ver**
    - **user\_rating**
    - **user\_rating\_ver**
    - **cont\_rating**
    - **prime\_genre** , transformo a numérico enumerándolo (el primer genero tendría el número 1, y así sucesivamente)
    - **sup\_devices.num**
    - **ipadSc\_urls.num**
    - **lang.num**
    - **vpp\_lic**
  2. **Búsqueda de vecinos más cercanos (KNN):** Es simplemente una variación de la búsqueda en profundidad de grafos, la distancia usada fue la euclidiana (comparando datos de una misma dimensión), se identificaba como un vecino cercano (Esto se repite hasta encontrar todos los vecinos pedidos) :
    - Si aún no se llena el conjunto de vecinos pedidos (10)
    - Si la distancia del nodo actual era más corta que la del peor nodo actual.
    - Luego se revisan los sub-árboles para evitar revisar todos los nodos.

A continuación se provee del código de ambas etapas.

```

1. def insertarRec(self,root,df_row,depth):genero
2.     global dimensions
3.     if(root == None):
4.         root = _node.create_node_charge_data(df_row,self.genre_dict,depth)
5.         return root
6.         cd = depth % dimensions
7.         df_app_ = self.fix_df(copy.deepcopy(df_row))
8.         if(self.special_case(df_app_.iloc[0][cd]) < self.special_case(root.df
_row.iloc[[0][0]][cd])):
9.             root.izq = self.insertarRec(root.izq,df_row,depth+1)
10.        else:
11.            root.der = self.insertarRec(root.der,df_row,depth+1)
12.        return root

```

Figura 1.1 inserción en K-d Tree

Aquí se da la cantidad de dimensiones, por medio de la variable global **dimensions**, y se usa la función **special\_case** , para corregir problemas con el tipo de dato.

```

1. def knn(self,root,dato,n):
2.     global dimensions
3.     dimension_ = None
4.     puntos_ordenados = list()
5.     pila = list()
6.     pila.insert(0,root)
7.     while(len(pila) > 0):
8.         nodo = pila.pop(0)
9.         dimension_ = nodo.depth % dimensions
10.        if((len(puntos_ordenados) < n)):
11.            puntos_ordenados.insert(0,nodo)
12.            self.ordenar(puntos_ordenados,dato.iloc[0][dimension_],dimension_)
13.        elif((self.euclidian_distance(dato.iloc[0][dimension_],nodo.df_row.ilo
c[[0][0]][dimension_]) < self.euclidian_distance(dato.iloc[0][dimension_],punt
os_ordenados[0].df_row.iloc[[0][0][dimension_]))) :
14.            puntos_ordenados.pop(0)
15.            puntos_ordenados.insert(0,nodo)
16.            self.ordenar(puntos_ordenados,dato.iloc[0][dimension_],dimension_)
17.        if(dato.iloc[0][dimension_] < nodo.df_row.iloc[[0][0]][dimension_]):
18.            pila.insert(0,nodo.der)
19.            pila.insert(0,nodo.izq)
20.        else:
21.            pila.insert(0,nodo.der)
22.            pila.insert(0,nodo.izq)
23.        pila = self.quitar_null(pila,root)
24.    return puntos_ordenados

```

Figura 2.1 Búsqueda de vecinos más cercanos

Igual que en la función pasada el número de dimensiones se da por la variable global, por otro lado para mantener el orden en la lista de los vecinos identificados más cercanos se ordena del peor al mejor (según distancia) por la función **ordenar**.