

# Proyecto Final

DOCENTE	CARRERA	CURSO
PhD(c) Vicente Machaca Arceda	Maestría en Ciencia de la Computación	Algoritmos y Estructura de Datos

## 1. Datos de los estudiantes

- Grupo: 07
- Integrantes:
  - Acuña Melo, Edgar
  - López Torres Herrera, Luis Rodrigo
  - Luna Flores, Julio Paolo
  - Vilela Arias, Roger

## 2. Competencias del proyecto

- Analizar e implementar algoritmos eficientes para la solución de problemas computacionales.
- Implementar estructuras de datos adecuadas, según el tipo de problema.

## 3. Equipos y materiales

- Javascript
- Navegador Web
- Cuenta en Github
- IDE de desarrollo

## 4. Trabajo de investigación

En el siguiente trabajo se realizó la implementación del clasificador KD-Tree teniendo como base a lo realizado en la practica 04 del curso, siendo el clasificador por frecuencia de palabras el seleccionado para el proyecto y desarrollado en su totalidad, orientandolo a la clasificacion de correos spam y no spam.

Para la compañía de antivirus ESET, el spam es cualquier forma de comunicación no solicitada que se envía de forma masiva (correo electrónico masivo no solicitado, o UBE). Su forma más frecuente es un correo electrónico de publicidad enviado a un gran número de direcciones (correo electrónico de publicidad no solicitado, o UCE), pero el "spamming" también existe a través de mensajes instantáneos, de texto (SMS), redes sociales o incluso mensajes de voz. Enviar spam es ilegal en la mayoría de jurisdicciones. Una de las vías más comunes para propagar contenido no solicitado es a través de botnets, grandes redes de dispositivos "zombieinfectados". A veces los correos en cadena y fraudulentos también se consideran spam, aunque difieren en que en general se reenvían por gente con buenas intenciones. (Fuente: [www.eset.com](http://www.eset.com))



Figura 01: Spam (Fuente: www.eset.com)

## 5. Desarrollo del proyecto

1. El proyecto final se encuentra desarrollado en el archivo TF-IDF.js el cual desarrolla el siguiente contenido.

```
1  let tokenize=(text)=>{
2    const punctuation = ['!', '#', '"', '%', '$', '"', '&', ')', '(', '+', '*', '-', ',', '/', '.
3    punctuation.forEach((e) => {
4      text = String(text.replaceAll(e, ""));
5    });
6    return text.toLowerCase().split(/\s+/g);
7  }
8  let dictionary=(tokens,dict)=>{
9    tokens.forEach((token)=>{
10     if(!dict.includes(token)){
11       dict.push(token);
12     }
13   })
14   return dict;
15 };
16 let vsm=(tokens,dict)=>dict.map((token)=>tokens.reduce((acc,curr)=>curr=== token?acc +1:acc,0));
17 let tf=(vsm,numberOfTokens)=>vsm.map((token)=>token/numberOfTokens);
18 let idf = (documentTokens,dict)=>dict.map((word)=>Math.log(documentTokens.length/documentTokens.
19 let tfidf=(tf,idf)=>tf.map((element,index)=>element*idf[index]);
20 let cosine=(tfIdf1,tfIdf2)=>tfIdf1.reduce((acc,curr,index)=>acc+curr*tfIdf2[index],0) / (Math.sq
21 let dict = [];
22 let review = [];
23 let adtfidf = [];
24 let atk = [];
```

```
25 function readFile(input) {
26   let file = input.files[0];
27   let reader = new FileReader();
28   reader.readAsText(file);
29   reader.onload = function() {
30     var content = reader.result;
31     var lines = content.split("\r\n");
32     for (var count = 1; count < lines.length; count++) {
33       const obj = {};
34       var rowContent = lines[count].split("|");
35       obj.text = String(rowContent[1]);
36       obj.classifier = String(rowContent[0]);
37       review.push(obj);
38     }
39     construccionVectores();
40     construccionKdTree();
41   };
42   reader.onerror = function() {
43     console.log(reader.error);
44   };
45 }

47 function construccionVectores() {
48   console.table("review");
49   atk = review.map((ele, index)=>{
50     const obj = [];
51     let dTokens = tokenize(ele.text);
52     dict=dictionary(dTokens,dict);
53     dTokens.forEach((e, i) => {
54       obj[i] = e;
55     });
56     return obj;
57   });
58
59   console.log(dict);
60
61   const adVms = [];
62   adtfidf = review.map((ele, index)=>{
63     const obj = [];
64     let dTokens = tokenize(ele.text);
65     let dVsm = vsm(dTokens,dict);
66     adVms[index] = dVsm;
67     let dtf = tf(dVsm,atk[index].length);
68     let didf = idf(atk,dict);
69     let dtfidf = tfidf(dtf,didf);
70     dtfidf.forEach((e, i) => {
71       obj[i] = e;
72     });
73     return obj;
74   });
```

```
76 console.table("adVms");
77 console.table("adtfidf");
78
79 review.map((ele, index) => {
80     ele.vector = adtfidf[index];
81 });
82 console.table("review");
83 }
84
85 function construccionKdTree() {
86     var oKdTree = new KdTree(adtfidf[0].length);
87     var root = oKdTree.build_kdtree_classifier(review);
88     console.log(root);
89
90     let querys = ["this is an amazing food, but is my personal opinion", "FREE Account... for three months. Dont miss out!"];
91     querys.forEach((query) => {
92         let queryTokens = tokenize(query);
93         let queryVsm = vsm(queryTokens, dict);
94         let querytf = tf(queryVsm, queryTokens.length);
95         let didf = idf(atk, dict);
96         let querytfidf = tfidf(querytf, didf);
97
98         console.log(querytfidf);
99
100        console.log("Consulta: ", query);
101        console.log("Resultado: ", oKdTree.classifierKNN(review, querytfidf, 3));
102    });
```

2. Así también se ha generado el archivo descriptor.py y se ha usado el archivo datamin.csv

```
1 import csv
2 import pandas as pd
3 import numpy as np
4 import string
5 from sklearn.feature_extraction.text import CountVectorizer
6 from sklearn.model_selection import train_test_split
7
8 spam_or_ham = pd.read_csv("data.csv", sep='|', encoding='latin-1')[["v1", "v2"]]
9 spam_or_ham.columns = ["label", "text"]
10 punctuation = set(string.punctuation)
11 print punctuation
12
13 ##Eliminamos los simbolos de puntuacion
14 def tokenize(sentence):
15     tokens = []
16     for token in sentence.split():
17         new_token = []
18         for character in token:
19             if character not in punctuation:
20                 new_token.append(character.lower())
21         if new_token:
22             tokens.append("".join(new_token))
23     return tokens
24
```

```

25 spam_or_ham.head()["text"].apply(tokenize)
26 train_text, test_text, train_labels, test_labels = train_test_split(spam_or_ham["text"],
27                                                                    spam_or_ham["label"],
28                                                                    stratify=spam_or_ham["label"])
29
30 print train_text
31 print "=====
32 print train_text.astype('U')
33
34 real_vectorizer = CountVectorizer(tokenizer = tokenize, binary=True)
35 train_X = real_vectorizer.fit_transform(train_text.astype('U'))
36 test_X = real_vectorizer.transform(test_text.astype('U'))
37
38 print(train_X)

```

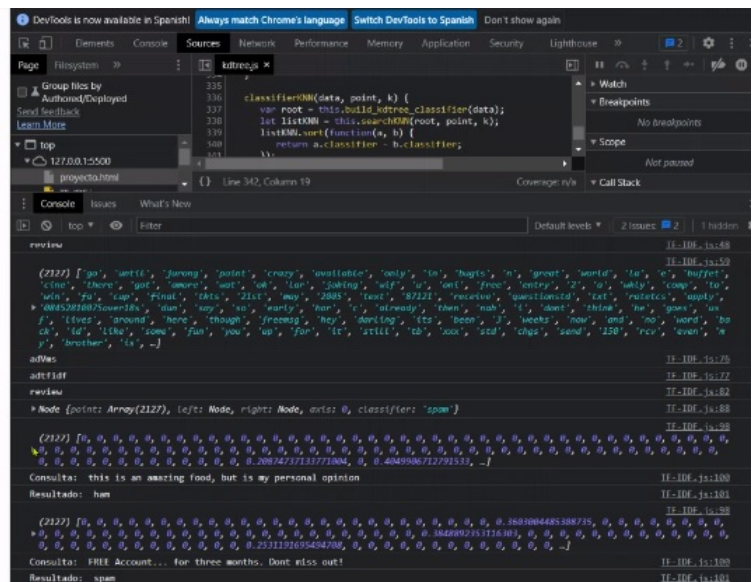
3. Cargando el archivo .csv se puede comprobar el funcionamiento adecuado del proyecto

### Procesamiento de Texto para diferenciar entre texto normal y SPAM

Seleccionar archivo data\_min.csv

### Procesamiento de Texto para diferenciar entre texto normal y SPAM

Seleccionar archivo data\_min.csv



## 6. Implementación

El desarrollo de la practica se utilizo el mismo repositorio compartido para la practica 04 debido a que se tomo cmo base lo desarrollado en dicha practica, los archivos de codigo fuente del desarrollo de la presente practica, se encuentran alojados en el siguiente enlace:

<https://github.com/JulioLunaUNSA/INFORME04.git>



📁 imagenes	PRUEBA VECTORIZACION IMAGENES
📄 Practica 04 - Grupo 7.pdf	Informe 04
📄 README.md	Initial commit
📄 TF-IDF.js	Ajustes TF-IDF
📄 data.csv	añadiendo descriptor y metodo de predicion
📄 data_min.csv	CLASIFICADOR NLP + AJUSTE KDTREE
📄 descriptor.py	CLASIFICADOR NLP + AJUSTE KDTREE
📄 kdtree.js	Ajustes TF-IDF
📄 main.html	Commit inicial
📄 opencv.js	PRUEBA VECTORIZACION IMAGENES
📄 p5.min.js	Commit inicial
📄 predicion.js	añadiendo descriptor y metodo de predicion
📄 proyecto.html	Ajustes TF-IDF
📄 prueba.js	Ajustes TF-IDF
📄 sketch.js	CLASIFICADOR NLP + AJUSTE KDTREE

## 7. Conclusiones

- De acuerdo a los objetivos planteados en el curso, con el desarrollo del proyecto final se logra un correcto analisis e implementacion de algoritmos con eficiencia demostrada para la solución de problemas computacionales.
- Se logro implementar estructuras de datos adecuadas manteniendo como base el KD-Tree, según el tipo de problema que se selecciono adicionalmente tomando como base las otras practicas desarrolladas.
- Se pudo discriminar los diferentes algoritmos teniendo en cuenta las prestaciones que se necesitan para la resolución de problemas computacionales.