



**UNIVERSIDAD NACIONAL
SAN AGUSTIN DE AREQUIPA**



Facultad de Ingeniería, Producción y Servicios

Escuela Profesional de Ciencia de la Computación

Laboratorio 1 - La Maldición de la Dimensionalidad

Presentado por:

Parizaca Mozo, Paul Antony

CUI:

20210686

Github:

<https://github.com/PaulParizacaMozo/Estructuras-De-Datos-Avanzadas>

Curso:

Estructuras de Datos Avanzadas

Docente:

Rosa Paccotacya Yanque

Arequipa, Perú

2023

Descripción del laboratorio

- Generar 100 puntos aleatorios entre 0 y 1 de dimensión d
- Calcular la distancia entre todos los pares de puntos (Distancia Euclidiana) (Hint 4950 distancias)
- Generar un histograma (pueden usar Python) de las distancias obtenidas para cada dimensión

Introducción

En el presente informe se va a generar 100 puntos aleatorios entre 0 y 1 en distintas dimensiones, específicamente en las siguientes dimensiones: 10, 50, 100, 500, 1000, 2000, 5000. Luego se va a calcular la distancia entre todos los pares ordenados generados y una vez obtengamos esas distancias generamos un histograma para analizar la tendencia de las distancias obtenidas.

Metodología

Para la generación de los puntos aleatorios se usará el lenguaje de programación c++, en el cual crearemos un programa para la generación de puntos, el programa consta de una clase Espacio la cual tiene como finalidad crear dos archivos de texto: data.txt y distancias.txt.

Una vez tengamos los dos archivos creados se ejecutará automáticamente un programa en python el cual usará como entrada los dos archivos para la generación de los histogramas haciendo uso de la librería Matplotlib.

El código de dichos programas se encuentra en el siguiente repositorio de Github:

<https://github.com/PaulParizacaMozo/Estructuras-De-Datos-Avanzadas>

Ejemplo de ejecución del programa:

-Ejecutamos el siguiente comando en consola donde se encuentren nuestros programas

```
> g++ -g -o eject main.cpp -pedantic && ./eject
```

-Nos pedirá ingresar la dimensión en la cual se quiere trabajar y la cantidad de puntos a generar

```

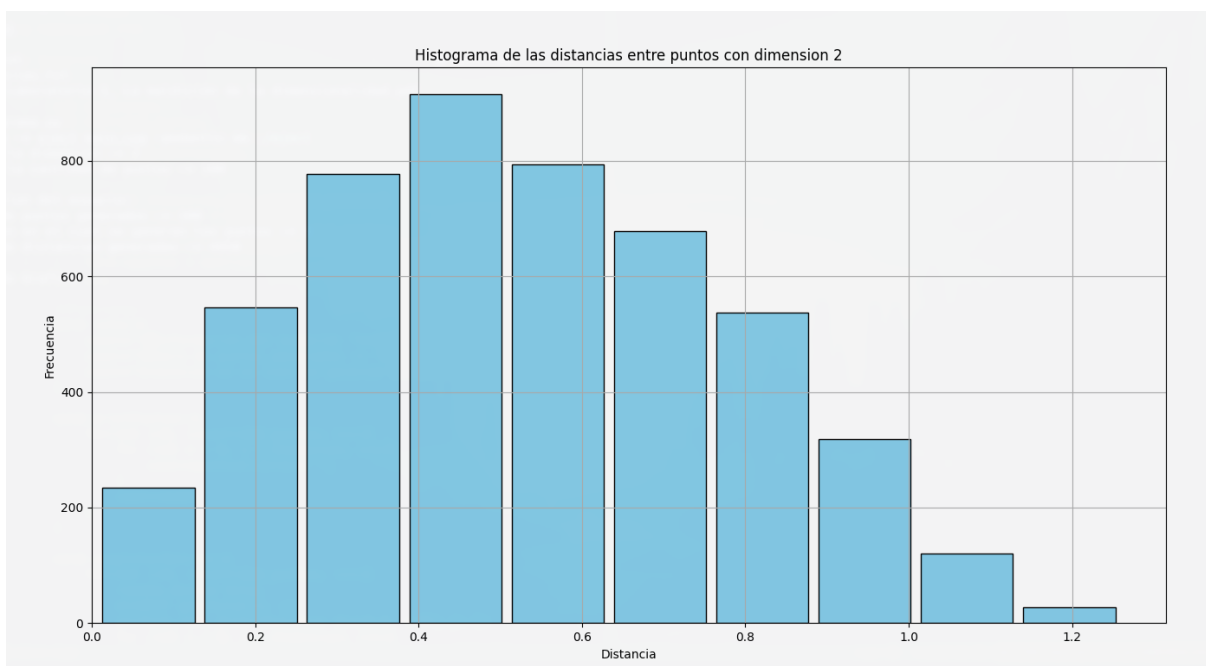
> g++ -g -o eject main.cpp -pedantic && ./eject
Ingrese la dimension -> 2
Ingrese la cantidad de puntos -> 100

Informacion del espacio:
Numero de puntos generados -> 100
Dimension en el cual se generan los puntos -> 2
Numero de distancias generadas -> 4950

Iniciando Grafica ...

```

-Nos mostrará la gráfica



-Y con eso culminará la ejecución del programa

```

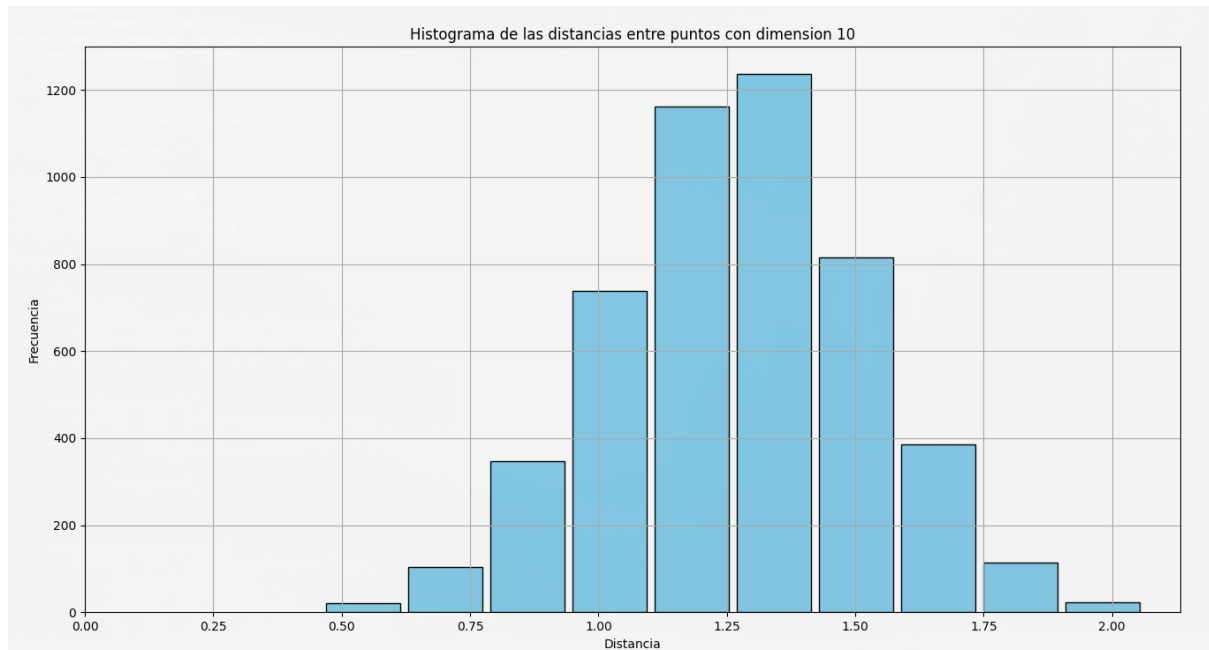
Finalizado.

^> ~/U/6semestre/E/Lab1 > on git P main !4 > took ⌚ 2m 1s |

```

Resultados

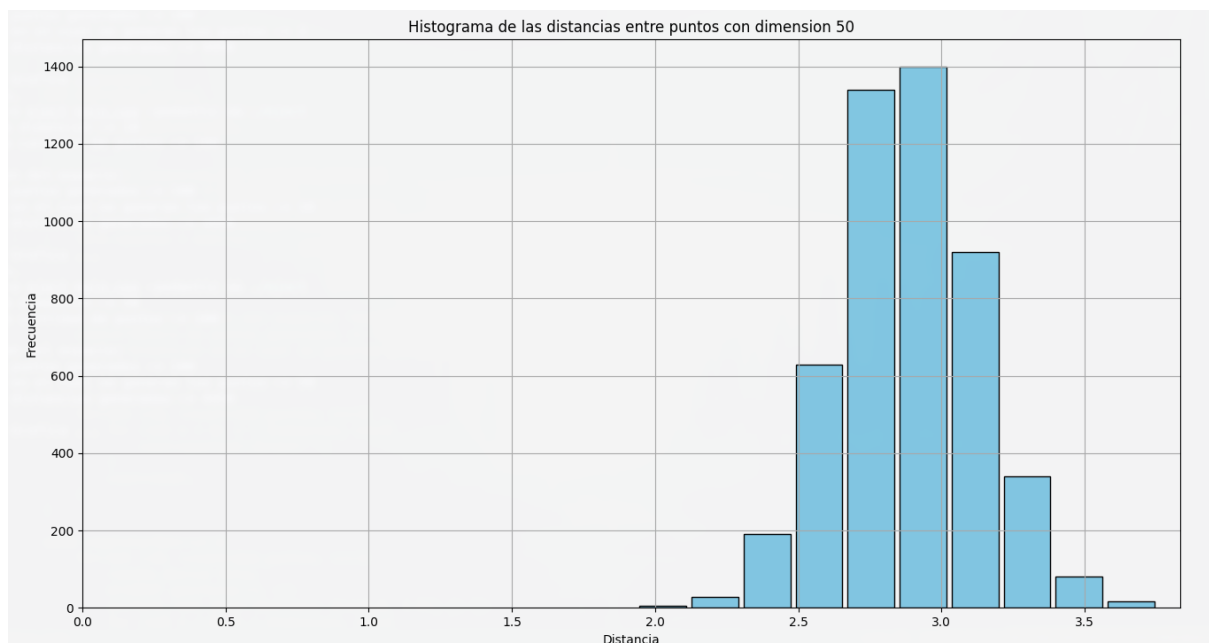
-Dimensión 10



Eje x: Los valores de las distancias se encuentran entre 0.25 y 2.25.

Eje y: Teniendo como frecuencia máxima un valor menor a 1300 y superior a 1200.

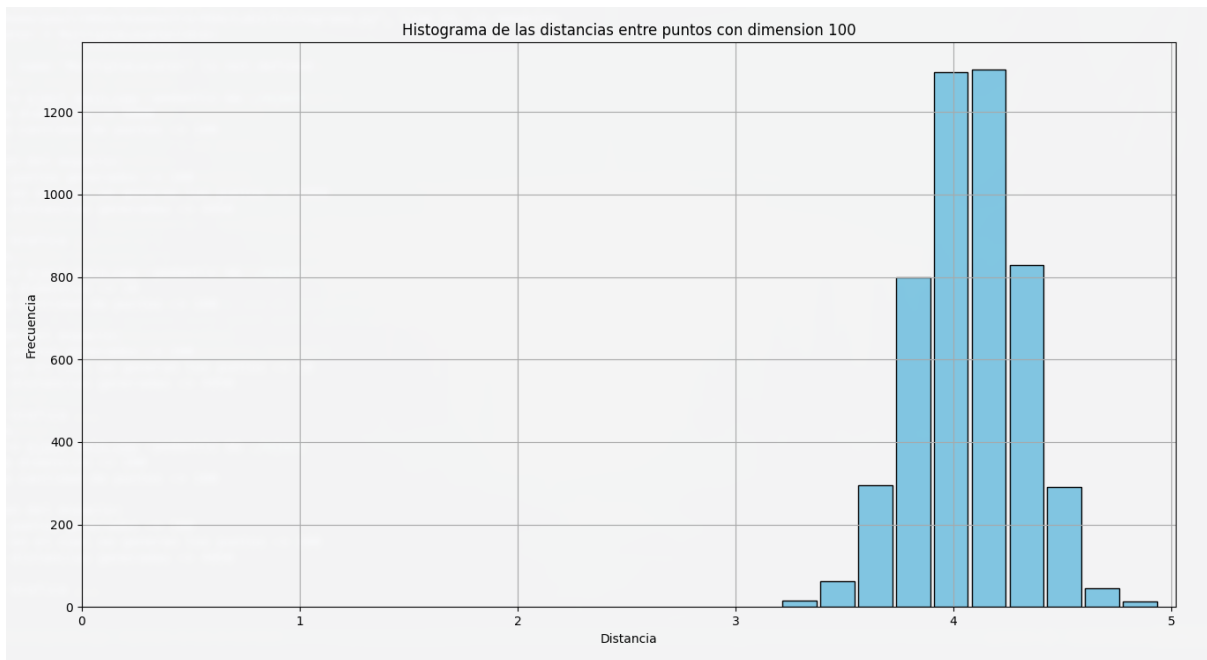
-Dimensión 50



Eje x: Los valores de las distancias se encuentran entre 1.5 y 4.0.

Eje y: Teniendo como frecuencia máxima un valor de 1400.

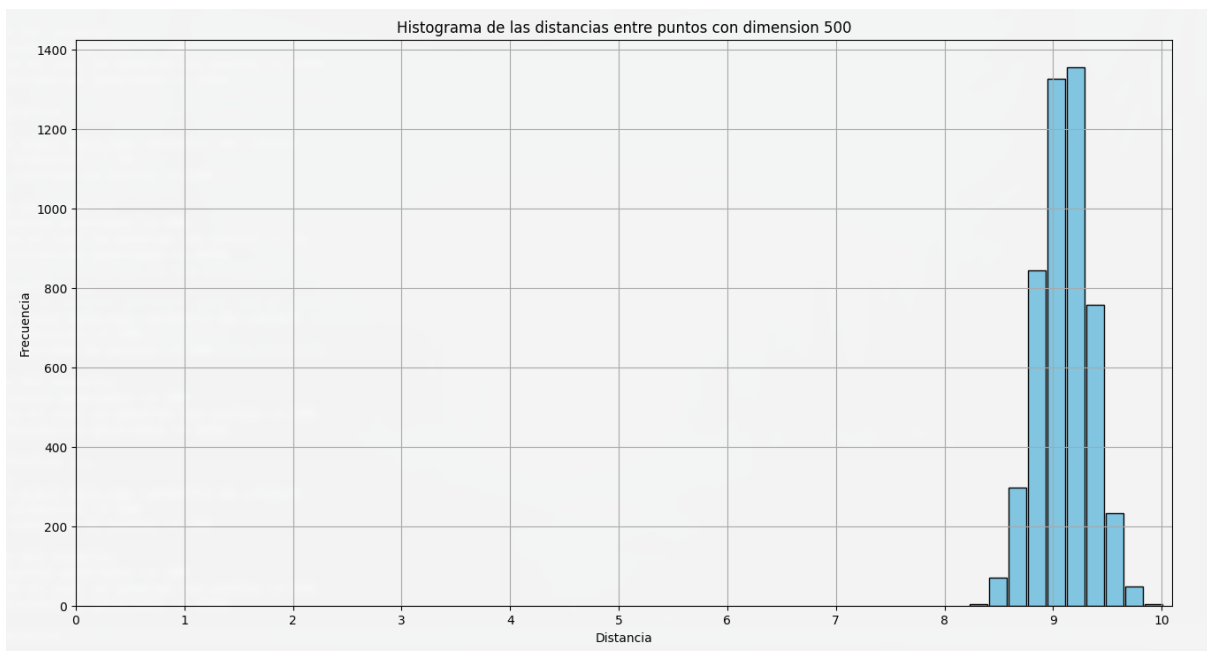
-Dimensión 100



Eje x: Los valores de las distancias se encuentran entre 3.0 y 5.0.

Eje y: Teniendo como frecuencia máxima un valor menor a 1300 y superior a 1200.

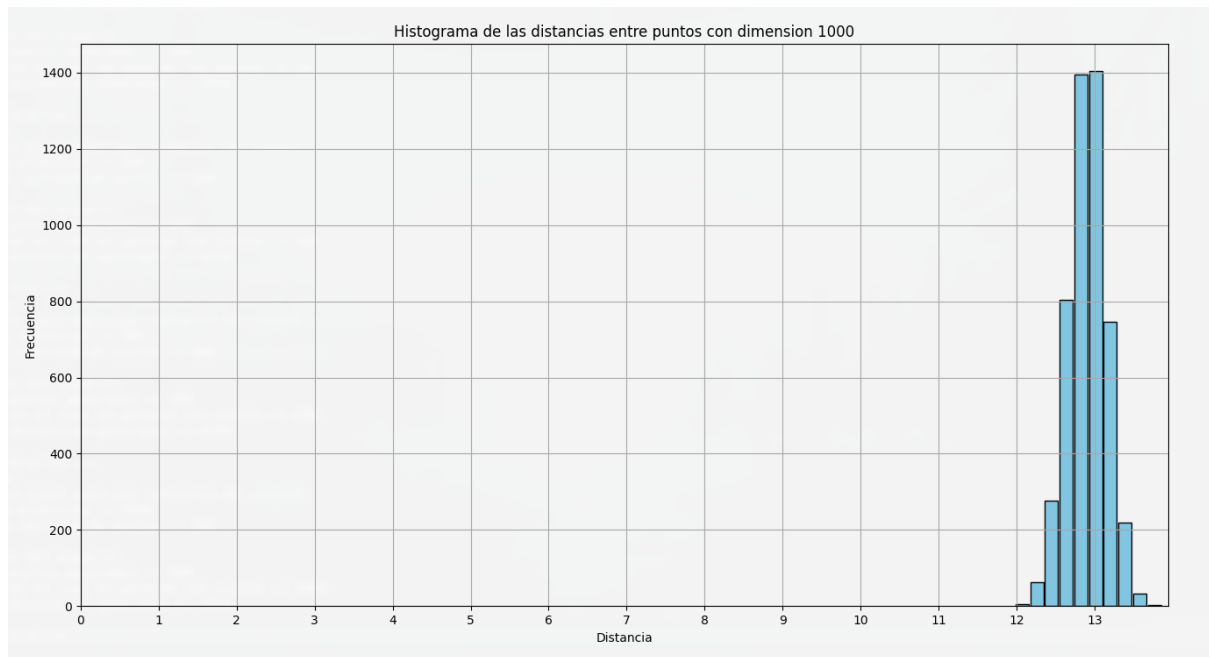
-Dimensión 500



Eje x: Los valores de las distancias se encuentran entre 8.0 y 10.0.

Eje y: Teniendo como frecuencia máxima un valor menor a 1400 y superior a 1300.

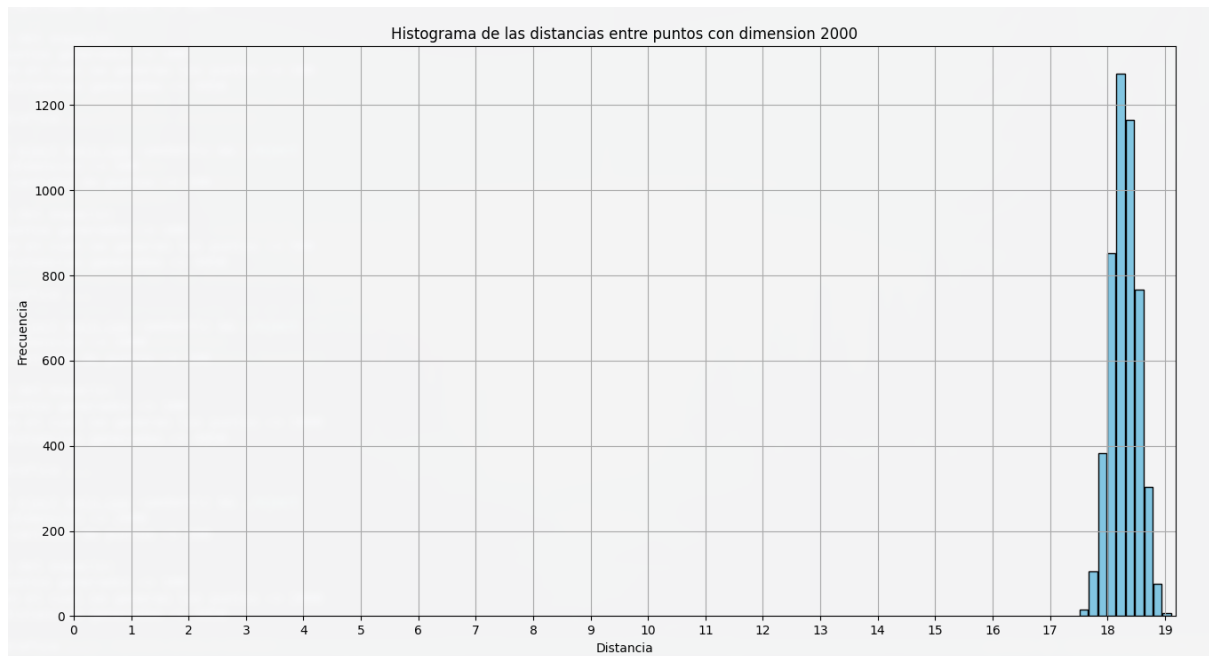
-Dimensión 1000



Eje x: Los valores de las distancias se encuentran entre 12.0 y 14.0.

Eje y: Teniendo como frecuencia máxima un valor que ronda los 1400.

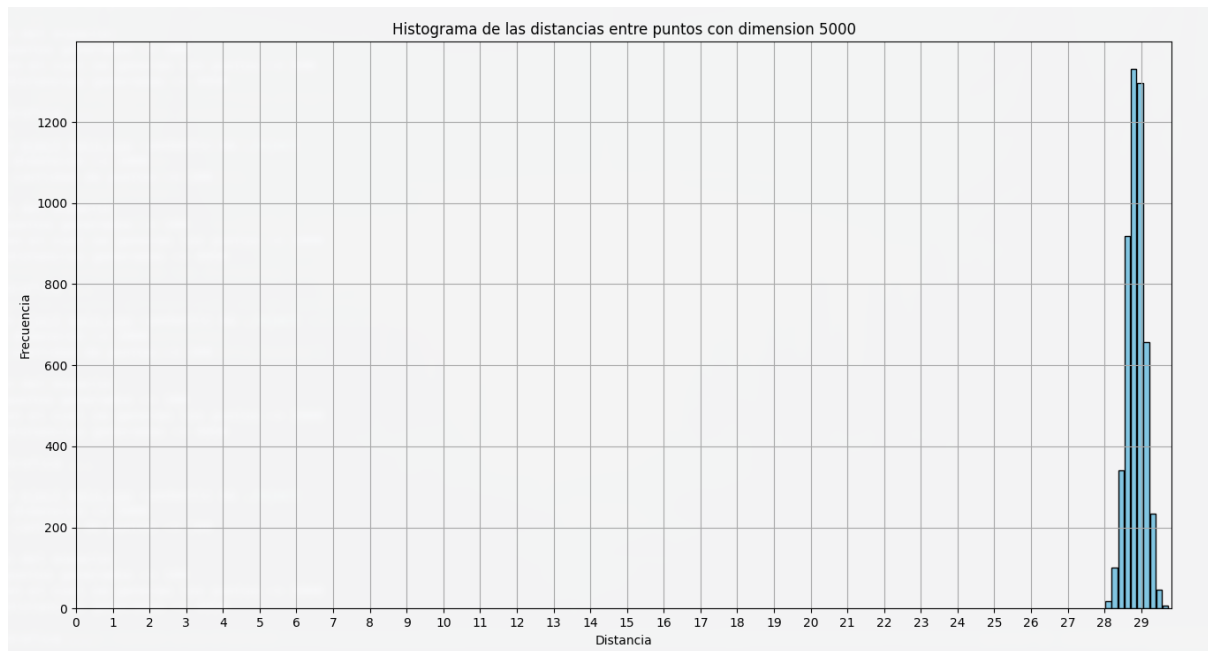
-Dimensión 2000



Eje x: Los valores de las distancias se encuentran entre 17.5 y 19.5.

Eje y: Teniendo como frecuencia máxima un valor menor a 1300 y superior a 1200.

-Dimensión 5000



Eje x: Los valores de las distancias se encuentran entre 28.0 y 30.0.

Eje y: Teniendo como frecuencia máxima un valor menor a 1300 y superior a 1200.

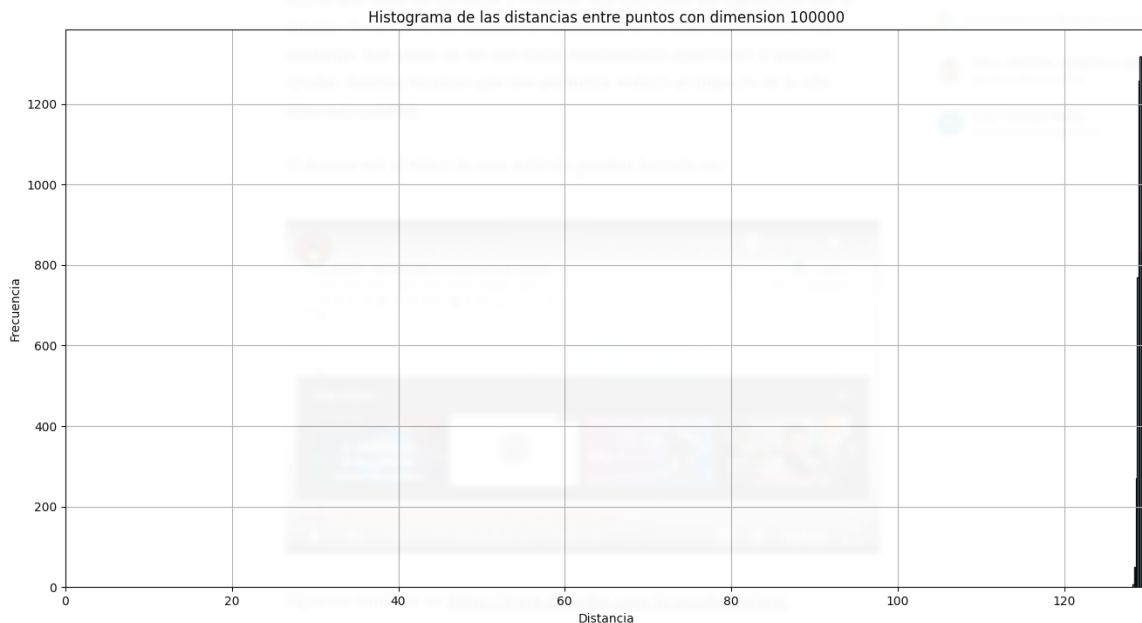
Análisis

Teniendo en cuenta los histogramas de dimensiones 10,50,100,500,1000,2000,5000 respectivamente podemos notar que hay un cambio significativo en las distancias de un punto con otro punto del mismo espacio que varía según se incrementa el número de dimensiones para el espacio.

Haciendo una comparación entre las distancias de dimensión 10 y las distancias de dimensión 1000 podemos ver que las distancias de dimensión 10 oscilan entre 0.25 y 2.25 mientras que las distancias de la dimensión 1000 oscilan entre 12.0 y 14.0. Con esto podemos ver que las distancias entre puntos van en incremento. Y comparando más aún las distancia de la dimensión 10 con las distancias de la dimensión 5000 en la cual las distancias entre un punto y otro oscilan entre 28.0 y 30.0 se nota el incremento significativo entre una distancia de una dimensión baja comparado con una dimensión alta.

Por ello también podemos decir que las distancias en las dimensiones bajas son representativas, que queremos decir con esto, que las distancias entre puntos pueden proporcionar información útil sobre cómo se relacionan los datos entre sí. Y a medida que se va creciendo el espacio en dimensiones las distancias entre los puntos se empieza a expandir, a crecer y por ello se van volviendo menos representativas y con un incremento

aún mayor de dimensiones se va perdiendo en su totalidad esta característica en la siguiente imagen podemos ver las distancias con una dimensión de 100000.



Viendo el histograma las distancias son muy similares entre ellas, se puede decir que entre un punto y otro las distancias son equidistantes.

Conclusión

- La alta dimensionalidad puede ser un problema y es necesario entender el efecto que tiene en los datos y como se ven afectados los algoritmos por lo mismo.[3]
- Las distancias entre puntos en dimensiones bajas son más representativas, significativas.
- A medida que las dimensiones incrementan en el espacio dejan de ser representativas, significativas.
- En un espacio con dimensiones altas tienden a ser similares, indistinguibles en algunos casos y con esto perder el cómo se relacionan las distancias entre puntos la una con la otra.

Código

Para el siguiente programa se trabajó con c++ y python para la elaboración del programa

–main.cpp

En el cual tenemos una clase Espacio


```

main.cpp
40 #include <random>
39 #include <vector>
38 #include <math.h>
37 #include <fstream>
36 using namespace std;
35
34 class Espacio{
33 public:
32     //Atributos
31     int dimension;
30     int cantidad_puntos;
29     vector<vector<double>>>puntos;
28     vector<double>distancias;
27
26     //Constructor
25     Espacio(int dimension,int cantidad_puntos){
24         this->dimension = dimension;
23         this->cantidad_puntos = cantidad_puntos;
22         vector<std::vector<double>>> matriz(cantidad_puntos,vector<double>(dimension)); n: value: n:
21         this->puntos = matriz;
20     }

```

La cual tiene atributos como dimensión, cantidad_puntos, un vector de vectores de tipo double puntos y un vector de tipo double para almacenar las distancias generadas.

La clase tiene funciones

info(): Muestra información del espacio(dimensión y cantidad de puntos)

```

15 //Funciones
14 void info(){
13     cout<<"\nInformacion del espacio: "<<endl;
12     cout<<"Numero de puntos generados -> "<<puntos.size()<<endl;
11     cout<<"Dimension en el cual se generan los puntos -> "<<puntos[0].size()<<endl;
10 }

```

completarMatriz(): Completa la matriz con números generados aleatoriamente entre 0 y 1

```

10 void completarMatriz(){
9     for(int i=0; i<cantidad_puntos; i++){
8         random_device hola;
7         mt19937 gen(hola()); sd:
6         uniform_real_distribution<> dis(0.0, 1.0); a: b:
5         for (int n = 0; n < dimension; ++n){
4             puntos[i][n] = dis(gen);
3         }
2     }
1 }

```

distancia euclidiana(): Algoritmo para calcular la distancia euclidiana

```

1 double distancia_euclidiana(vector<double>punto1,vector<double>punto2){
2     double distancia = 0;
3     for (int i=0 ;i<dimension;i++) {
4         distancia += pow(punto2[i]-punto1[i],2); x: y:
5     }
6     distancia = sqrt(distancia); x:
7     return distancia;
8 }

```

calcular_distancias(): Calcula todas las distancias de todos los puntos con los demas puntos

```
10 void calcular_distancias(){
11     for(int i=0 ;i<cantidad_puntos;i++){
12         for(int j=0;j<cantidad_puntos;j++){
13             if(i < j){
14                 double aux = distancia_euclediana(puntos[i],puntos[j]); punto1: punto2:
15                 distancias.push_back(aux); x:
16             }
17         }
18     }
19 }
```

escribir_archivos(): Escribe dos archivos: data.txt y distancias.txt que serán usados por el programa escrito en python para la generación de histogramas.

```
12 void escribir_archivos(){
11     ofstream file;
10     file.open("distancias.txt"); s:
9     for(auto vec : distancias){ : double
8         file<<vec<<'\n';
7     }
6     file.close();
5     file.open("data.txt"); s:
4     file<<dimension<<'\n';
3     file<<cantidad_puntos<<'\n';
2     file.close();
1     }
84 };
```

main(): Solicita al usuario ingresar la dimensión y la cantidad de puntos para el espacio y con esa información crear un objeto de la clase espacio y llamar a las funciones necesarias para mostrar el histograma.

```
18 int main()
17 {
16     int dim,cant;
15     cout<<"Ingrese la dimension -> ";
14     cin>>dim;
13     cout<<"Ingrese la cantidad de puntos -> ";
12     cin>>cant;
11     Espacio a(dim,cant); dimension: cantidad_puntos:
10     a.info();
9     a.completarMatriz();
8     a.calcular_distancias();
7     cout<<"Numero de distancias generadas -> "<<a.distancias.size()<<endl<<endl;
6     cout<<"Iniciando Grafica ..."<<endl;
5     a.escribir_archivos();
4     cin.get();
3     system("python histograma.py"); command:
2     cout<<"Finalizado."<<endl;
1     return 0;
105 }
```

-Histograma.py

```
histograma.py
1 import matplotlib.pyplot as plt
2 from matplotlib.ticker import MultipleLocator
3
4 dist = []
5 data = []
6
7 with open('distancias.txt', 'r') as archivo:
8     for linea in archivo:
9         valor_distancia = float(linea.strip())
10        dist.append(valor_distancia)
11
12 with open('data.txt', 'r') as archivo:
13     for linea in archivo:
14         valor = int(linea.strip())
15         data.append(valor)
16
17 plt.hist(dist, edgecolor='k', rwidth=0.9, color='skyblue')
18 plt.xlabel('Distancia')
19 plt.ylabel('Frecuencia')
20 plt.title('Histograma de las distancias entre puntos con dimension '+str(data[0]))
21
22 dim = 1
23 x_locator = MultipleLocator(dim)
24 plt.gca().xaxis.set_major_locator(x_locator)
25
26 plt.xlim(0)
27 plt.grid(True)
28 plt.show()
```

El programa en python leerá la información de los archivos data.txt y distancias.txt y los almacenará en vectores para el manejo de la información que luego será requerida para la creación de los histogramas. Se usó la librería matplotlib para la generación de los histogramas. Código Completo en el repositorio de Github:

<https://github.com/PaulParizacaMozo/Estructuras-De-Datos-Avanzadas>

Referencias

1.- uniform_real_distribution

https://en.cppreference.com/w/cpp/numeric/random/uniform_real_distribution

2.- Distancia Euclidiana.

<https://www.lifeder.com/distancia-euclidiana/>

3.- La Maldición de la Dimensionalidad

<https://medium.com/@nicolasarrioja/la-maldici%C3%B3n-de-la-dimensionalidad-f7a6248cf9a>

[a](#)

4.- La maldición de la dimensionalidad - 3 - Ciencia de Datos y Machine Learning

<https://youtu.be/WlavTxJSDUs>