



CIENCIA DE LA COMPUTACIÓN

ESTRUCTURAS DE DATOS AVANZADA

---

## ORIENTACIÓN DE PUNTOS - PCA

---

Docente Carlos Eduardo Atencio Torres

Jun 2023

Autor:  
Iris Rocio Curo Quispe

Semestre VI  
2023-1

# Orientación de puntos - PCA

Iris Rocio Curo Quispe  
iris.curo@ucsp.edu.pe

June 18, 2023

## Abstract

En el presente trabajo se presenta la implementación del cálculo de PCA y el vector director de un conjunto de puntos. Para graficar los puntos se usó Python de forma que se pueda visualizar la nube de puntos junto con el vector director hallado para estos puntos, el cual será el primer eigenvector.

**Keywords**— PCA, vector director, números aleatorios.

## 1 Introducción

En el presente informe se describe la implementación de los códigos proporcionados en el repositorio de Github de *Diego Mazala* para realizar el Análisis de Componentes Principales (PCA, por sus siglas en inglés) utilizando el lenguaje de programación C++ y la librería Eigen3. El objetivo de esta implementación es calcular los eigenvectores de una nube de puntos tridimensionales y visualizarlos en un gráfico 3D.

El PCA es una técnica fundamental en el campo del procesamiento de datos y la reducción de dimensionalidad. Permite encontrar las direcciones principales en un conjunto de datos y explicar la selección de los mismos en términos de estas direcciones. En este informe, se presenta un análisis detallado del código `pca_test.cpp`, el cual implementa el PCA utilizando la biblioteca Eigen en C++. El objetivo principal de este informe es comprender cómo se realiza el cálculo del PCA, cómo se generan los resultados y qué información se puede obtener a partir de ellos.

El código `pca_test.cpp` comienza solicitando al usuario que especifique el tamaño de la matriz de datos. A continuación, genera una matriz de datos aleatorios utilizando una distribución uniforme. Los valores de esta matriz se guardan en un archivo de texto para su posterior análisis. Luego, el código utiliza la biblioteca Eigen para realizar el cálculo del PCA. Se emplea la clase `pca_t` proporcionada, que encapsula las operaciones necesarias para obtener los resultados del PCA. Una vez realizado el cálculo, los resultados se guardan en el mismo archivo de texto previamente creado.

En el informe, se tomarán y analizarán los resultados obtenidos a través del cálculo del PCA. Se mostrarán matrices como la matriz de entrada, la matriz centrada, la matriz de covarianza y la matriz de proyección. Además, se examinarán los vectores propios y los valores propios asociados al PCA, los cuales son fundamentales para comprender las direcciones principales en los datos y su contribución a la necesidad del conjunto. Además, se realizará una reproyección de la matriz de datos y se calculará el error entre la matriz original y la reproyección obtenida. Este error nos brinda información sobre la calidad de la representación de los datos utilizando el PCA y nos permite evaluar la efectividad de la técnica.

Para lograr esto, se realizaron los siguientes pasos:

- Generación de puntos aleatorios y vector de dirección: Se agregó una funcionalidad adicional para generar puntos aleatorios y un vector de dirección usando el primer *eigenvector* en C++. Se implementó un código que genera estos datos de manera aleatoria junto con el vector director calculado y los guarda en un archivo para su posterior utilización.
- Modificación del código para leer los datos generados: Por último, se modificó el código original en Python para leer los datos generados por el código en C++. Esto implicó realizar cambios en la forma en que se cargan y procesan los datos en el código Python, de manera que se puedan utilizar los datos generados en C++ para realizar la gráfica que los presenta.

## 2 Generación de puntos aleatorios y vector director

A continuación se presenta una explicación del código `pca_test.cpp` en cuestión. El código es un ejemplo de cómo realizar el PCA utilizando la biblioteca Eigen en C++:

- El programa comienza verificando si se ha proporcionado un argumento en la línea de comandos, que representa el tamaño de la matriz de datos a generar. Si no se proporciona el argumento adecuado, se muestra un mensaje de uso y se finaliza el programa.
- A continuación, se genera una matriz de datos aleatorios llamada `pca_data_matrix` utilizando la biblioteca Eigen. La matriz tiene "m" filas y 3 columnas (3D), donde "m" es el tamaño especificado por el usuario. Los valores de la matriz se generan utilizando una distribución uniforme en el rango de -5.0 a 5.0, este rango se puede modificar.
- Después de generar la matriz de datos, se guarda en un archivo de texto llamado `pca_data.txt`. El archivo contiene la matriz de datos y se estructura de la siguiente manera: los valores de la matriz se escriben fila por fila, separados por espacios. A continuación, se realiza el cálculo del PCA utilizando la clase `pca_t` la cual esta definida en `pca.h`.
- El resultado del PCA se guarda en el archivo `pca_data.txt`. Primero, se obtienen los vectores propios (*eigenvectors*) de la matriz de covarianza y se guarda el primer vector propio en el archivo. Luego, se cierra el archivo y se muestra un mensaje indicando que la matriz y el vector director han sido guardados exitosamente en el archivo `pca_data.txt`.
- A continuación, se realiza el cálculo del PCA nuevamente utilizando la clase `pca_t` y se muestran los resultados en la consola. Se imprimen las siguientes matrices y vectores: la matriz de entrada, la matriz centrada, la matriz de covarianza, la matriz de proyección, la matriz de medias, los valores propios (*eigenvalues*) y los vectores propios (*eigenvectors*). Estos resultados proporcionan información detallada sobre los diferentes pasos y resultados del PCA.
- Además, se realiza una reproyección de la matriz de datos utilizando la función `reprojection()` y se calcula el error entre la matriz original y la reproyección utilizando la función `norm()`. El resultado de la reproyección y el error se imprimen en la consola.

Listing 1: Código `pca_test.cpp` que genera una matriz de datos aleatorios, realiza el cálculo del PCA utilizando la biblioteca Eigen en C++.

```
1 #include <iostream>
2 #include <fstream>
3 #include <random>
4 #include "pca.h"
5
6 #include <experimental/filesystem>
7
8 using namespace std;
9 namespace fs = std::experimental::filesystem;
10
11 int main(int argc, char* argv[])
12 {
13
14     if (argc != 2)
15     {
16         std::cout << "Usage: " << argv[0] << "(matrix size) m x 3 e.g. 10 x 3"
17             << std::endl;
18         return EXIT_FAILURE;
19     }
20
21     int m = std::stoi(argv[1]);
22
23     // Generar puntos aleatorios en pca_data_matrix
24     std::random_device rd;
25     std::mt19937 gen(rd());
26     std::uniform_real_distribution<float> dis(-5.0, 5.0);
27
28     Eigen::Matrix<float, Eigen::Dynamic, 3> pca_data_matrix(m, 3);
```

```

28     for (int i = 0; i < pca_data_matrix.rows(); ++i)
29     {
30         for (int j = 0; j < pca_data_matrix.cols(); ++j)
31         {
32             pca_data_matrix(i, j) = dis(gen);
33         }
34     }
35
36     // Guardar la matriz en un archivo de texto
37     std::ofstream outfile("pca_data.txt");
38     if (outfile.is_open())
39     {
40         outfile << "Puntos:\n";
41         for (int i = 0; i < pca_data_matrix.rows(); ++i)
42         {
43             for (int j = 0; j < pca_data_matrix.cols(); ++j)
44             {
45                 outfile << pca_data_matrix(i, j) << " ";
46             }
47             outfile << std::endl;
48         }
49
50         pca_t<float> pca;
51         pca.set_input(pca_data_matrix);
52         pca.compute();
53
54         Eigen::Matrix<float, Eigen::Dynamic, Eigen::Dynamic> eigen_vectors = pca
55             .get_eigen_vectors();
56         Eigen::Matrix<float, Eigen::Dynamic, 1> first_eigenvector =
57             eigen_vectors.col(0);
58
59         // Escribir el primer eigenvector en el archivo
60         outfile << "Vector Director:\n";
61         for (int i = 0; i < first_eigenvector.rows(); ++i)
62         {
63             outfile << first_eigenvector(i) << " ";
64         }
65         outfile << std::endl;
66
67         outfile.close();
68         std::cout << "Matriz y Vector Director guardados en pca_data.txt" << std
69             ::endl;
70     }
71     else
72     {
73         std::cerr << "No se pudo abrir el archivo pca_data.txt" << std::endl;
74         return EXIT_FAILURE;
75     }
76
77     // Realizar el calculo de PCA
78     pca_t<float> pca;
79     pca.set_input(pca_data_matrix);
80     pca.compute();
81
82     std::cout
83     << "Input Matrix:      \n" << pca.get_input_matrix() << std::endl << std
84         ::endl
85     << "Centered Matrix:   \n" << pca.get_centered_matrix() << std::endl <<
86         std::endl
87     << "Covariance Matrix: \n" << pca.get_covariance_matrix() << std::endl
88         << std::endl
89     << "Projection Matrix: \n" << pca.get_projection_matrix() << std::endl
90         << std::endl
91     << "Mean Matrix:      \n" << pca.get_mean() << std::endl << std::endl
92     << "Eigen Values:     \n" << pca.get_eigen_values() << std::endl << std
93         ::endl

```

```

86         << "Eigen Vectors:      \n" << pca.get_eigen_vectors() << std::endl <<
            std::endl;
87
88     const auto& reprojection = pca.reprojection();
89     auto error = (pca_data_matrix - reprojection).norm();
90
91
92     std::cout
93         << "Reprojected Matrix:\n" << reprojection << std::endl << std::endl
94         << std::fixed
95         << "Error:                \n" << error << std::endl << std::endl;
96
97     return EXIT_SUCCESS;
98 }

1 Eigen::MatrixXf generarFigura(int n)
2 {
3     std::random_device rd;
4     std::mt19937 gen(rd());
5     std::uniform_real_distribution<float> dis(0, 2 * M_PI);
6
7     Eigen::MatrixXf figura(n, 3);
8     for (int i = 0; i < n; ++i)
9     {
10         float x = static_cast<float>(i - n / 2);
11         float theta1 = dis(gen);
12         float theta2 = dis(gen);
13         float y = x * std::sin(theta1);
14         float z = x * std::sin(theta2);
15
16         figura(i, 0) = x;
17         figura(i, 1) = y;
18         figura(i, 2) = z;
19     }
20
21     return figura;
22 }

```

Listing 2: Función para generar una distribución tipo mariposa de puntos aleatorios

### 3 Generación de gráfico 3D

El código presentado en `pca_example.py` tiene como objetivo cargar los datos de un archivo llamado `tpca_data.txt` para graficar los puntos y un vector director en un espacio tridimensional, como ejemplo, se puede ver la imagen 2. A continuación, se presenta una descripción paso a paso del funcionamiento del código:

- En primer lugar, se cargan los datos de un archivo llamado `pca_data.txt`, que contiene la información de los puntos y el vector director. Los puntos se almacenan en un arreglo bidimensional llamado *"points"*, mientras que el vector director se guarda en un arreglo unidimensional llamado *"direction"*.
- A continuación, se calculan los puntos extremos de la nube de puntos utilizando las funciones `np.min()` y `np.max()`. Estos puntos extremos se utilizarán para determinar la longitud y posición del vector director en el gráfico.
- La longitud del vector director se calcula utilizando la función `np.linalg.norm()`, que devuelve la norma euclidiana del arreglo. Luego, se determina un factor de escala dividiendo la norma de la distancia entre los puntos extremos de la nube de puntos por la longitud del vector director. Este factor de escala se reafirma para ajustar la longitud del vector director proporcionalmente a la distancia entre los puntos extremos.

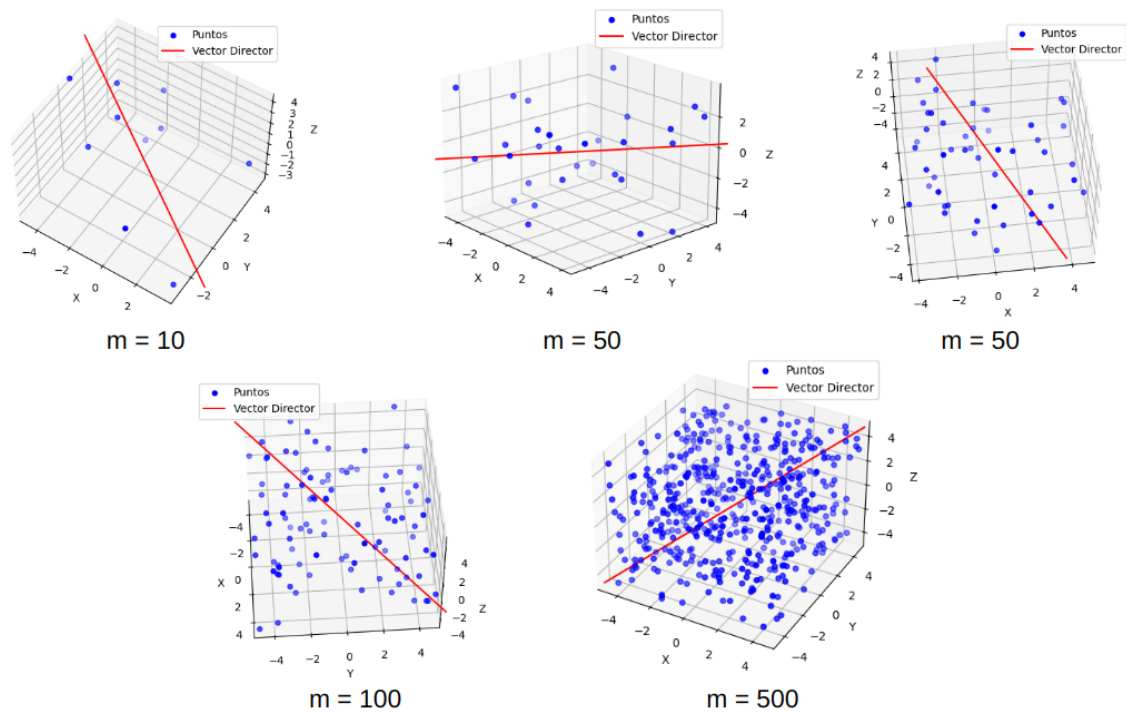


Figure 1: Gráfico que muestra  $m$  puntos aleatorios generados junto con el vector director calculado usando el primer *eigenvector* al aplicar PCA.

- A continuación, se calculan los puntos de inicio y fin del vector director. El punto de inicio se determina como el punto medio entre los puntos extremos de la nube de puntos, mientras que el punto de fin se obtiene sumando el vector director escalado al punto de inicio.
- Luego, se crea una figura y unos ejes 3D utilizando la biblioteca matplotlib. Los puntos se grafican como una dispersión tridimensional utilizando la función `scatter()`, y se les asigna el color azul.
- Finalmente, el vector director se dibuja en el gráfico usando la función `plot3D()`. Los puntos de inicio y fin del vector se sustentan como argumentos para definir las coordenadas  $X$ ,  $Y$  y  $Z$  de la línea que representa el vector director. Esta línea se dibuja en color rojo.
- Se configuran los límites de los ejes del gráfico utilizando los puntos extremos de la nube de puntos, lo que garantiza que todos los puntos y el vector director sean visibles en el gráfico.
- Se agregan etiquetas a los ejes  $X$ ,  $Y$  y  $Z$ , y se agrega una leyenda que indica la representación de los puntos y el vector director en el gráfico.

Finalmente, se muestra el gráfico resultante utilizando la función `show()`.

Listing 3: Implementación de `pca_example.py` para graficar la nube de puntos con su vector director.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from mpl_toolkits.mplot3d import Axes3D
4
5 # Cargar los datos desde el archivo pca_data.txt
6 with open('pca_data.txt', 'r') as file:
7     lines = file.readlines()
8
9 # Extraer los puntos y el vector director
10 points = np.loadtxt(lines[1:-2])
11 direction = np.fromstring(lines[-1], dtype=float, sep=' ')
12
13 # Calcular los puntos extremos de la nube de puntos
14 min_point = np.min(points, axis=0)
15 max_point = np.max(points, axis=0)
16

```

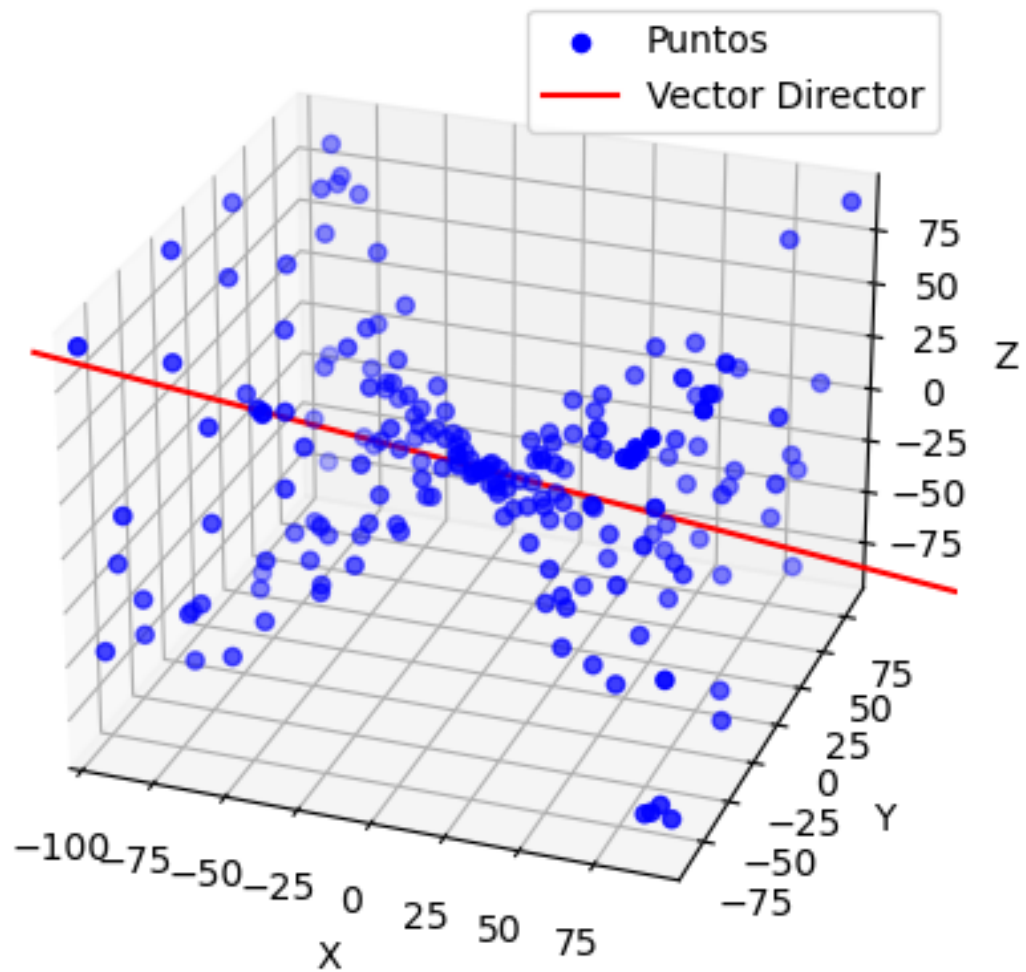


Figure 2: Gráfico que muestra  $m = 200$  puntos aleatorios generados con la distribución mariposa junto con el vector director calculado usando el primer *eigenvector* al aplicar PCA.

```

17 # Calcular la longitud del vector director
18 point_range = max_point - min_point
19 direction_length = np.linalg.norm(direction)
20 scaling_factor = np.linalg.norm(point_range) / direction_length
21 scaled_direction = direction * scaling_factor
22
23 # Calcular los puntos de inicio y fin del vector director
24 start_point = (min_point + max_point) / 2.0 - scaled_direction / 2.0
25 end_point = start_point + scaled_direction
26
27 # Crear la figura y los ejes 3D
28 fig = plt.figure()
29 ax = fig.add_subplot(111, projection='3d')
30
31 # Graficar los puntos
32 ax.scatter(points[:, 0], points[:, 1], points[:, 2], c='blue', label='Puntos')
33
34 # Graficar el vector director proporcionalmente
35 ax.plot3D([start_point[0], end_point[0]],
36           [start_point[1], end_point[1]],
37           [start_point[2], end_point[2]], 'red', label='Vector Director')
38

```

```

39 # Configurar los l mites de los ejes
40 ax.set_xlim(min_point[0], max_point[0])
41 ax.set_ylim(min_point[1], max_point[1])
42 ax.set_zlim(min_point[2], max_point[2])
43
44 # Configurar etiquetas y leyendas
45 ax.set_xlabel('X')
46 ax.set_ylabel('Y')
47 ax.set_zlabel('Z')
48 ax.legend()
49
50 # Mostrar el gr fico
51 plt.show()

```

```

(base) triscuro@triscuro-spin-5P314-S1:~/Documentos/2023/IDA/PCAV1$ g++ -std=c++11 -I /usr/include/eigen3 pca_test.cpp -o pca_test
(base) triscuro@triscuro-spin-5P314-S1:~/Documentos/2023/IDA/PCAV1$ ./pca_test
Matriz y Vector Director guardados en pca_data.txt
Input Matrix:
-4.07457  3.33403  -1.58595
-4.38237  0.864254  0.175915
 3.33276  3.11406  -1.65239
-0.663545  4.37966  -3.74565
 1.35425  1.83809  1.28752
 2.3817  2.71948  -0.915281
 0.827351  -1.82179  -4.44267
-3.64792  2.46844  -0.855593
 1.24965  1.29245  1.99082
 3.82782  4.94681  -4.55854

Centered Matrix:
-4.88728  1.02128  -0.8837767
-4.39588  -1.44849  1.59809
 3.32804  0.801313  -0.238221
-0.676256  2.06691  -2.32347
 1.34154  -0.474657  2.70969
 2.29099  0.40673  0.506973
 0.814639  -4.13454  -3.02049
-3.66063  0.147687  0.566581
 1.23693  -1.0203  3.41299
 3.8151  2.63486  -3.13636

Covariance Matrix:
 9.41109  0.958671  -1.48404
 0.958671  3.72685  -0.848764
-1.48404  -0.848764  5.17129

Projection Matrix:
-3.59883  -4.86062  3.31324  0.471396  0.318009  2.0517  0.950632  -3.56316  -0.0985022  5.01613
 1.88153  -0.309747  -0.715367  2.96474  -2.99835  -1.12253  0.933225  0.866634  -3.75797  2.25783
-1.12438  0.497173  -0.478004  -1.05644  -0.525563  -0.448935  5.01067  -0.543954  -0.295793  -1.03477

Mean Matrix:
0.0127115  2.31275  -1.42217
0.0127115  2.31275  -1.42217
0.0127118  2.31275  -1.42217
0.0127117  2.31275  -1.42217
0.0127118  2.31275  -1.42217
0.0127118  2.31275  -1.42217
0.0127117  2.31275  -1.42217
0.0127118  2.31275  -1.42217
0.0127118  2.31275  -1.42217
0.0127118  2.31275  -1.42217

Eigen Values:
10.0954  4.8922  3.32163

Eigen Vectors:
 0.932368  -0.357706  0.0523123
 0.181982  0.339373  -0.922881
-0.312366  -0.869984  -0.381517

Reprojected Matrix:
-2.6491  5.2168  -2.91789
-4.73884  3.51377  -1.58826
 3.121  1.30866  -0.486286
 1.32175  4.06937  -3.73057
-0.0722634  1.63067  1.56209
 1.06161  1.58845  -0.107611
-0.496285  -2.06979  -4.14536
-2.98184  4.35466  -2.20084
-0.670615  1.32996  2.15368
 5.4237  2.18493  -2.84869

Error:
7.453445

```

Figure 3: Captura de pantalla que muestra los resultados luego de ejecutar el código `pca_test.cpp` con 10 puntos.



## 4 Conclusiones

Gracias al análisis del código `pca_test.cpp` me permitió comprender cómo se implementa y se utiliza el PCA en el lenguaje de programación C++ utilizando la biblioteca Eigen, además de su visualización usando Python. Los resultados obtenidos brindaron información valiosa sobre las direcciones principales en los datos mediante la técnica del PCA. Durante el análisis de los resultados, se evaluaron diferentes matrices, como la matriz de entrada, la matriz centrada, la matriz de covarianza y la matriz de proyección. Estas matrices nos proporcionaron información sobre la estructura y la adquisición de los datos originales.

Los vectores propios y los valores propios asociados al PCA. Estos vectores propios representan las direcciones principales en los datos, y los valores propios indican la contribución relativa de cada dirección principal a la necesidad total del conjunto de datos. Estos resultados son fundamentales para comprender y visualizar las principales características de los datos. El cálculo del error entre la matriz original y la matriz reproyectada permitió evaluar la calidad de la aproximación obtenida mediante el PCA. Un error bajo indica una buena aproximación de los datos originales, mientras que un error alto puede sugerir que las direcciones principales seleccionadas no son suficientes para representar adecuadamente la alteración de los datos.