

# Desarrollo Taller 1 AI: Decisión Bayesiana

Juan Sebastián Bravo Santacruz

## I. OBJETIVOS

- 1) Visualizar e interpretar datos de acuerdo con la estadística descriptiva de los mismos.
- 2) Diseñar e implementar clasificadores Bayesianos.

## II. CONJUNTO DE DATOS DATA\_2D

El conjunto de datos data\_2D (ver archivo data.npy) contiene datos pertenecientes a dos clases {a, b}.

A. *Grafique los datos utilizando un color distintivo para cada clase:*

Inicialmente se realizó la importación de los datos del archivo data.npy y se separaron los conjuntos data\_2D y data\_3D en 2 distintas variables:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn import model_selection
4
5 # Importar la informacion del archivo data
6 data = np.load('data.npy', allow_pickle=True)
7 data = data.item()
8
9 '''----- Punto 1 -----'''
10
11 # Separar los conjuntos de datos data_2D y data_3D
12 data_2D = data['data_2D']
13 data_3D = data['data_3D']
```

Una vez separados los conjuntos de datos, se extrajeron los valores de las clases a y b y se realizó la gráfica de estas dos clases usando el módulo matplotlib.

```
14
15 # Adquirir informacion de las 2 clases del conjunto de datos data_2D
16 data_2D_a = data_2D['data_a']
17 data_2D_b = data_2D['data_b']
18
19 # Grafica de los datos con colores distintivos para cada clase
20 plt.scatter(data_2D_a[0], data_2D_a[1], c='red', label='Clase a')
21 plt.scatter(data_2D_b[0], data_2D_b[1], c='blue', label='Clase b')
22 plt.title('Conjunto de datos data_2D ')
23 plt.xlabel('Atributo 1')
24 plt.ylabel('Atributo 2')
25 plt.legend()
26 plt.grid()
```

En este caso, se eligió el color rojo para la clase a y el color azul para la clase b para poder distinguirlos adecuadamente en la gráfica.

La gráfica obtenida se presenta en la siguiente imagen:

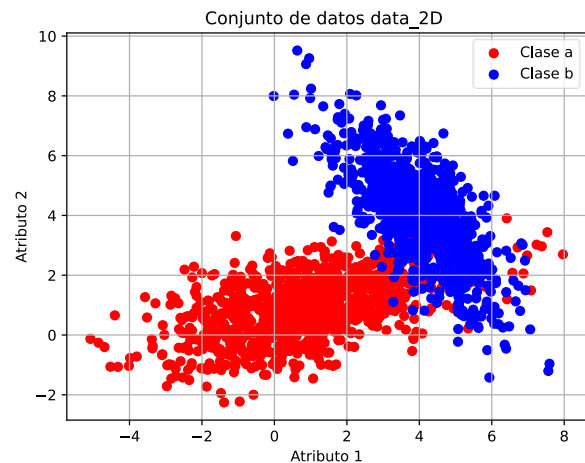


Figura 1. Conjunto de datos data\_2D.

B. *Determine el centro de cada clase  $\{\mu_a, \mu_b\}$ :*

Para encontrar el centro de cada una de las clases, se utilizó la función mean() del módulo numpy:

```
28
29 # Calculo del centro de cada clase
30 u_2D_a = data_2D_a.mean(axis=1)
31 u_2D_b = data_2D_b.mean(axis=1)
```

En este caso, al tener 2 atributos, se obtiene un vector de 2 valores de media, uno por cada atributo:

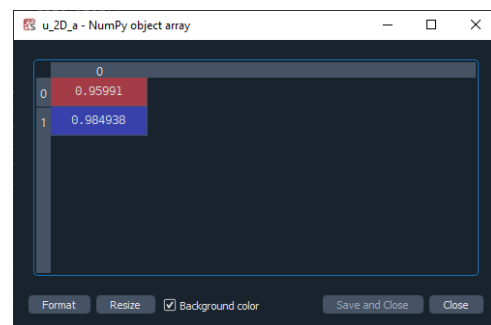


Figura 2. Vector medias para clase a.

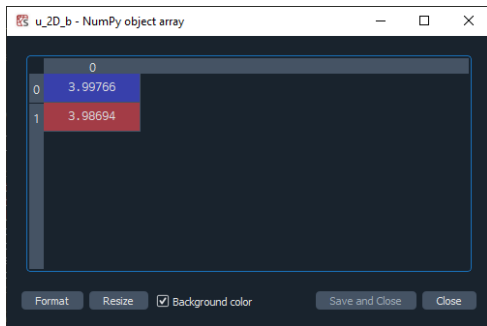


Figura 3. Vector medias para clase b.

C. Determine las matrices de covarianza de cada clase  $\{K_a, K_b\}$ . ¿Qué se puede concluir?:

Para encontrar las matrices de covarianza de cada una de las clases, se utilizó la función `cov()` del módulo `numpy`:

```
33 # Calculo de las matrices de covarianza de cada clase
34 K_2D_a = np.cov(data_2D_a)
35 K_2D_b = np.cov(data_2D_b)
```

En este caso, al tener 2 atributos, se obtiene una matriz cuadrada de 2x2:

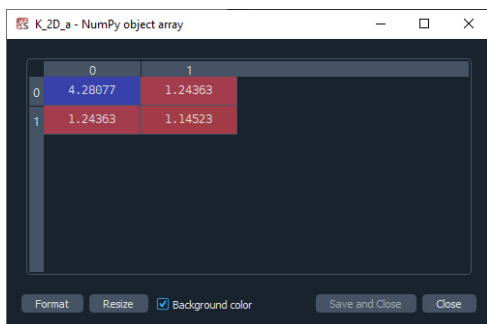


Figura 4. Matriz de covarianza clase a.

De la matriz de covarianza de la clase a, se puede concluir que los atributos no son independientes debido a que no es una matriz diagonal, y que tienen una relación directa (a medida que el atributo 1 aumenta, el atributo 2 aumenta), adicionalmente, se puede concluir que los datos varían más en el atributo 1 que en el atributo 2.

Lo anteriormente dicho, se puede comprobar al ver la gráfica de datos realizados en el inciso A. del ejercicio.

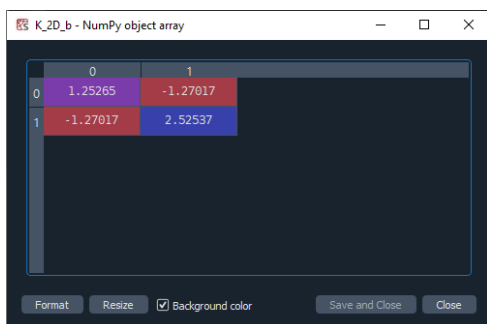


Figura 5. Matriz de covarianza clase b.

De la matriz de covarianza de la clase b, se puede concluir que los atributos no son independientes debido a que no es una matriz diagonal, y que tienen una relación inversa (a medida que el atributo 1 aumenta, el atributo 2 disminuye), adicionalmente, se puede concluir que los datos varían más en el atributo 2 que en el atributo 1.

Lo anteriormente dicho, se puede comprobar al ver la gráfica de datos realizados en el inciso A. del ejercicio.

D. Determine y visualice el histograma de los datos.

Se utilizó el módulo `matplotlib` para encontrar el histograma de los datos, haciendo uso de la función `histogram2d()` y `meshgrid` haciendo un histograma para cada una de las clases como se muestra en el siguiente código:

```
37 # Histograma clase a data_2D
38 fig = plt.figure()
39 ax = fig.add_subplot(projection='3d')
40 hist, xedges, yedges = np.histogram2d(data_2D_a[0], data_2D_a[1], bins=(20,20))
41 xpos, ypos = np.meshgrid(xedges[:1], yedges[:1])
42
43 xpos = xpos.flatten()/2.
44 ypos = ypos.flatten()/2.
45 zpos = np.zeros_like(xpos)
46
47 dx = xedges[1] - xedges[0]
48 dy = yedges[1] - yedges[0]
49 dz = hist.flatten()
50
51 ax.bar3d(xpos, ypos, zpos, dx, dy, dz, zsort='average')
52
53 plt.title('Histograma clase_a data_2D')
54 plt.xlabel('Atributo 1')
55 plt.ylabel('Atributo 2')
56
57 # Histograma clase b data_2D
58 fig = plt.figure()
59 ax = fig.add_subplot(projection='3d')
60 hist, xedges, yedges = np.histogram2d(data_2D_b[0], data_2D_b[1], bins=(20,20))
61 xpos, ypos = np.meshgrid(xedges[:1], yedges[:1])
62
63 xpos = xpos.flatten()/2.
64 ypos = ypos.flatten()/2.
65 zpos = np.zeros_like(xpos)
66
67 dx = xedges[1] - xedges[0]
68 dy = yedges[1] - yedges[0]
69 dz = hist.flatten()
70
71 ax.bar3d(xpos, ypos, zpos, dx, dy, dz, zsort='average')
72
73 plt.title('Histograma clase_b data_2D')
74 plt.xlabel('Atributo 1')
75 plt.ylabel('Atributo 2')
76 plt.show()
77
```

Los histogramas obtenidos se presentan en las siguientes imágenes:

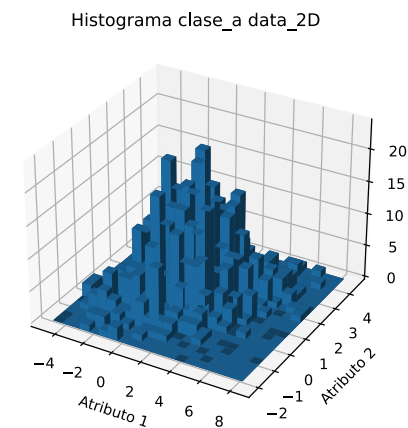


Figura 6. Histograma clase a data\_2D.

En el histograma realizado para la clase a, se puede observar que la mayor frecuencia de datos se encuentra entorno al valor central de los datos muy cercano a (1,1).

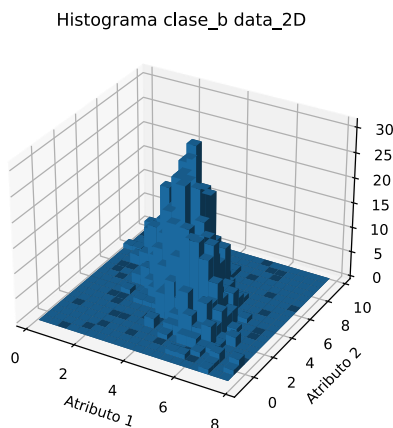


Figura 7. Histograma clase b data\_2D.

En el histograma realizado para la clase b, se puede observar que la mayor frecuencia de datos se encuentra entorno al valor central de los datos muy cercano a (4,4).

### III. CONJUNTO DE DATOS DATA\_3D

**A.** Divida los datos aleatoriamente en conjunto de prueba (20 %) y de entrenamiento (80 %).

La división de datos se realizó adquiriendo inicialmente los datos de cada una de las clases para el conjunto data\_3D, seguidamente se utilizó la función `train_test_split()` del módulo `sklearn` como se muestra en la imagen:

```
81 # Visualizar los datos de cada una de las clases del conjunto de datos data_3D
82 data_3D_a = data_3D['data_a']
83 data_3D_b = data_3D['data_b']
84
85 # División de datos en conjuntos de prueba y entrenamiento
86 training_a, testing_a = model_selection.train_test_split(data_3D_a, test_size=(0.2*len(data_3D_a)), train_size=(0.8*len(data_3D_a)))
87 training_b, testing_b = model_selection.train_test_split(data_3D_b, test_size=(0.2*len(data_3D_b)), train_size=(0.8*len(data_3D_b)))
```

En este, para el parámetro `test_size`, se encontró el 20% del número total de datos de cada clase y para el parámetro `train_size` se encontró el 80% del número total de datos de cada clase.

**B.** Visualice el conjunto de entrenamiento con un color para cada clase.

Para graficar el conjunto de entrenamiento de cada clase, se utilizó la función `scatter3D` pasando como parámetro las columnas de las variables de entrenamiento creadas para cada clase como se muestra en el siguiente código:

```
87 # Visualización de datos
88 ax = plt.axes(projection='3d')
89 ax.scatter3D(training_a[:,0], training_a[:,1], training_a[:,2], c='red', label='Clase a')
90 ax.scatter3D(training_b[:,0], training_b[:,1], training_b[:,2], c='blue', label='Clase b')
91 plt.title('Conjunto de entrenamiento datos data_3D')
92 ax.set_xlabel('Atributo 1')
93 ax.set_ylabel('Atributo 2')
94 ax.set_zlabel('Atributo 3')
95 plt.legend()
96 plt.grid()
97 plt.show()
```

La gráfica obtenida es la siguiente:

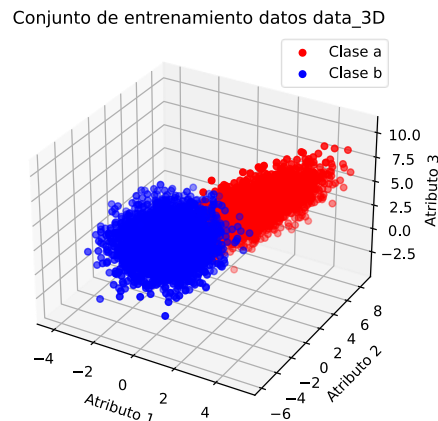


Figura 8. Conjunto de entrenamiento data\_3D.

**C.** Implemente un clasificador Bayesiano Gaussiano, y:

i. Estime la función de verosimilitud de cada clase:

Inicialmente para la implantación de la función de verosimilitud presentada en la siguiente ecuación:

$$f_X(X|c_i)P(c_i) = \frac{P(c_i)}{\sqrt{2\pi^n|K|}} e^{\left[-\frac{1}{2}(X-\mu)^T K^{-1}(X-\mu)\right]}$$

Se encontró el número de atributos y las probabilidades a priori de cada una de las clases, en este caso, asumiendo que estaban distribuidas homogéneamente:

```
99 # Numero de atributos
100 n = np.size(data_3D_a,axis=1)
101
102 # Probabilidades a priori
103 P_3D_a = len(training_a)/(len(training_a) + len(training_b))
104 P_3D_b = len(training_b)/(len(training_a) + len(training_b))
```

Posteriormente, se encontraron los valores centrales y las matrices de covarianza para cada una de las clases.

```
106 # Calculo del centro de cada clase
107 u_3D_a = training_a.mean(axis=0)
108 u_3D_b = training_b.mean(axis=0)
109
110 # Matriz de covarianza:
111 K_3D_a = np.cov(np.transpose(training_a))
112 K_3D_b = np.cov(np.transpose(training_b))
```

Al tener 3 atributos se tiene un vector con 3 valores de media, uno por cada atributo.

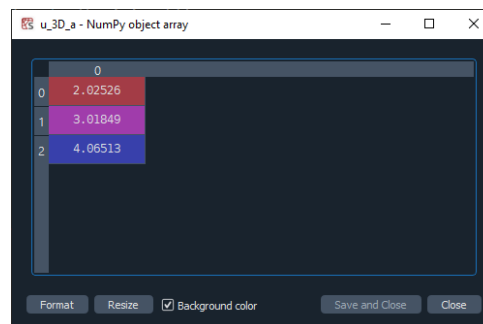


Figura 9. Vector medias para clase a.

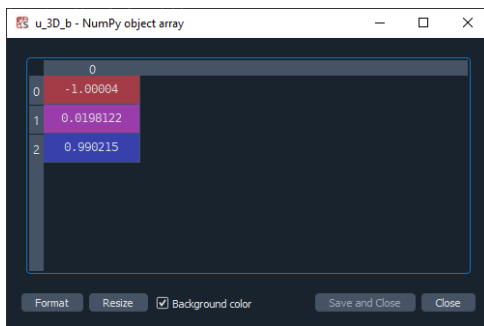


Figura 10. Vector medias para clase a.

Adicionalmente, se tienen matrices de covarianza de 3x3, debido al número de atributos:

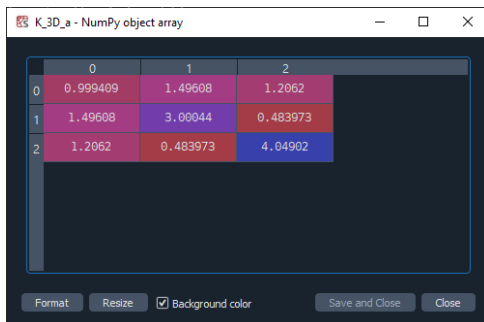


Figura 11. Matriz de covarianza clase a.

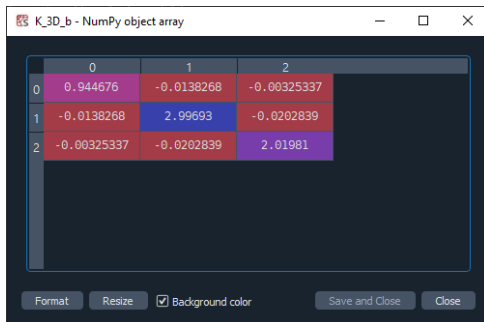


Figura 12. Matriz de covarianza clase b.

De las matrices de covarianza, se puede observar que en caso de la clase a, se tiene una dependencia entre atributos, mientras que para el caso de la clase b, se podría decir existe una dependencia mínima entre estos.

Una vez encontradas cada una de las constantes necesarias para encontrar la función, se creó una matriz de prueba que contenía los conjuntos de prueba de la clase a y b, se creó un vector con las etiquetas de clasificación reales para posteriormente calcular el error de la clasificación y finalmente un vector para almacenar la salida de la clasificación a realizar:

```
114 # Matriz de prueba
115 X = np.concatenate((testing_a, testing_b), axis=0)
116
117 # vector de comparacion
118 y_real = np.concatenate((np.zeros((len(testing_a),1)), np.ones((len(testing_b),1))), axis=0)
119
120 # vector con las clasificaciones
121 y_bayes = np.zeros((len(X),1))
```

Finalmente, se implementó la función, asignando un 0 para la clase a y un 1 para la clase b como se muestra en el siguiente código:

```
123 for i in range(len(X)):
124     # Funcion de verosimilitud de cada clase:
125     P_a_X = P_3D_a*(1/(np.sqrt(2*(np.pi**n)*np.linalg.det(K_3D_a))))*(
126         np.exp(-0.5*np.matmul(np.matmul(X[i,:], u_3D_a), np.linalg.inv(K_3D_a)), X[i,:], u_3D_a)))
127
128     P_b_X = P_3D_b*(1/(np.sqrt(2*(np.pi**n)*np.linalg.det(K_3D_b))))*(
129         np.exp(-0.5*np.matmul(np.matmul(X[i,:], u_3D_b), np.linalg.inv(K_3D_b)), X[i,:], u_3D_b)))
130
131     # Vector de probabilidades
132     P = (P_a_X, P_b_X)
133
134     # Proceso de clasificacion
135     clase = P.index(max(P))
136
137     # Clase a
138     if clase==0:
139         y_bayes[i] = 0
140
141     # Clase b
142     if clase==1:
143         y_bayes[i] = 1
144
```

- ii. Visualice la clasificación realizada sobre el conjunto de prueba:

Para la visualización del conjunto de prueba, inicialmente, se crearon 2 nuevas matrices para cada una de las clases en base a la clasificación realizada y seguidamente se graficó el conjunto haciendo uso de la función scatter3D():

```
145 # Clasificacion del conjunto de prueba
146 a_bayes = np.array([X[i] for i in range(len(y_bayes)) if y_bayes[i] == 0])
147 b_bayes = np.array([X[i] for i in range(len(y_bayes)) if y_bayes[i] == 1])
148
149 # Visualizacion conjunto de prueba clasificado
150 plt.figure()
151 ax = plt.axes(projection='3d')
152 ax.scatter3D(a_bayes[:,0], a_bayes[:,1], a_bayes[:,2], c='red', label='Clase a')
153 ax.scatter3D(b_bayes[:,0], b_bayes[:,1], b_bayes[:,2], c='blue', label='Clase b')
154 plt.title('Conjunto de prueba datos data_3D clasificado con Bayes Gaussiano')
155 ax.set_xlabel('Atributo 1')
156 ax.set_ylabel('Atributo 2')
157 ax.set_zlabel('Atributo 3')
158 plt.legend()
159 plt.grid()
```

La gráfica obtenida es la siguiente:

Conjunto de prueba datos data\_3D clasificado con Bayes Gaussiano

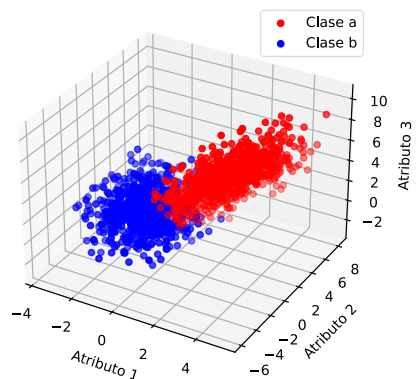


Figura 13. Conjunto de prueba clasificado con Bayes Gaussiano.

- iii. Determine el error de clasificación sobre el conjunto de prueba:

Se determinó el error de clasificación, encontrando la cantidad de clasificaciones mal realizadas dividido el número de elementos clasificados.

```
161 # Calculo error porcentual de la clasificacion
162 error_bayes = 100*sum(y_bayes != y_real)/len(y_real)
163 print(f"Error Bayes Gaussiano: {error_bayes[0]:.3f}%")
164
```

Los errores porcentuales obtenidos en 3 ejecuciones del programa fueron:

```
Error Bayes Gaussiano: 1.875%
Error Bayes Gaussiano: 2.875%
Error Bayes Gaussiano: 3.125%
```

#### D. Implemente un clasificador Bayesiano Gaussiano Naive, y:

- Estime la función de verosimilitud de cada clase:

En este caso, la función implementada en se presenta en la siguiente ecuación:

$$f_X(X|c_i)P(c_i) = \frac{P(c_i)^n}{\sqrt{2\pi^n(\sigma_1^2 \dots \sigma_n^2)}} e^{\left[-\frac{1}{2}\sum_{i=1}^n \left(\frac{x_i - \mu_i}{\sigma_i}\right)^2\right]}$$

Se encontró el número de atributos y las probabilidades a priori de cada una de las clases, en este caso, asumiendo que estaban distribuidas homogéneamente. Estas son las mismas que en el caso anterior. Adicionalmente, se encontraron las desviaciones estándar para cada una de las clases.

```
167 # Probabilidades a priori y medias iguales al clasificador anterior
168
169 # Desviaciones estandar
170 std_a = training_a.std(axis=0)
171 std_b = training_b.std(axis=0)
172
```

Al tener 3 atributos, se obtiene 1 desviación estándar por cada uno de los atributos:

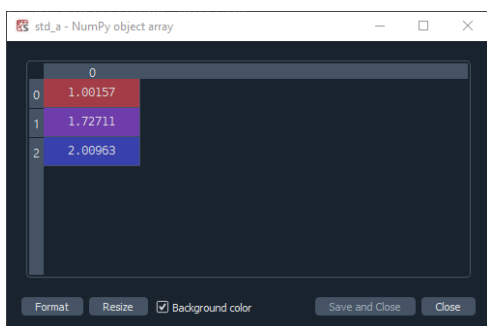


Figura 14. Desviaciones estandar clase a.

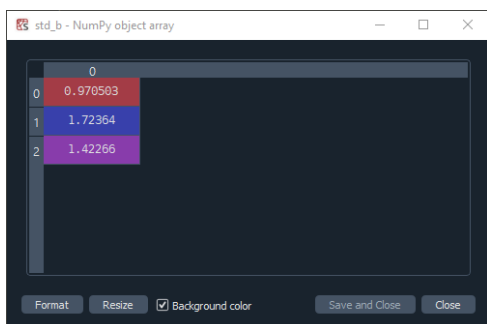


Figura 15. Desviaciones estandar clase b.

Se creó un vector con las etiquetas de clasificación reales para posteriormente calcular el error de la clasificación y finalmente un vector para almacenar la salida de la clasificación a realizar:

```
173 # vector de comparacion
174 y_real = np.concatenate((np.zeros((len(testing_a),1)), np.ones((len(testing_b),1))), axis=0)
175
176 # vector con las clasificaciones
177 y_bayes_naive = np.zeros((len(X),1))
```

Finalmente, se implementó la función, asignando un 0 para la clase a y un 1 para la clase b como se muestra en el siguiente código:

```
179 for i in range(len(X)):
180     # lista con probabilidades de atributos
181     p_a_X = []
182     p_b_X = []
183     for j in range(n):
184         # Calculo de probabilidades a posteriori de cada atributo
185         p_a_X.append(P_3D_a*(1/(np.sqrt(2*(np.pi**n))*std_a[j]))*(np.exp(-0.5*(X[i,j]-u_3D_a)**2/std_a[j]**2)))
186         p_b_X.append(P_3D_b*(1/(np.sqrt(2*(np.pi**n))*std_b[j]))*(np.exp(-0.5*(X[i,j]-u_3D_b)**2/std_b[j]**2)))
187
188     # Probabilidad por clase
189     p_a_X = np.prod(p_a_X)
190     p_b_X = np.prod(p_b_X)
191
192     # Vector de probabilidades
193     P = (p_a_X, p_b_X)
194
195     # Proceso de clasificación
196     clase = P.index(max(P))
197
198     # Clase a
199     if clase==0:
200         y_bayes_naive[i] = 0
201
202     # Clase b
203     if clase==1:
204         y_bayes_naive[i] = 1
205
```

- Visualice la clasificación realizada sobre el conjunto de prueba:

Para la visualización del conjunto de prueba, inicialmente, se crearon 2 nuevas matrices para cada una de las clases en base a la clasificación realizada y seguidamente se graficó el conjunto haciendo uso de la función scatter3D():

```
206 # Clasificación del conjunto de prueba
207 a_bayes_naive = np.array([X[i] for i in range(len(y_bayes_naive)) if y_bayes_naive[i] == 0])
208 b_bayes_naive = np.array([X[i] for i in range(len(y_bayes_naive)) if y_bayes_naive[i] == 1])
209
210 # Visualización conjunto de prueba clasificado
211 plt.figure()
212 ax = plt.axes(projection='3d')
213 ax.scatter3D(a_bayes_naive[:,0], a_bayes_naive[:,1], a_bayes_naive[:,2], c='red', label='Clase a')
214 ax.scatter3D(b_bayes_naive[:,0], b_bayes_naive[:,1], b_bayes_naive[:,2], c='blue', label='Clase b')
215 plt.title('Conjunto de prueba datos data_3D clasificado con Bayes Gaussiano Naive')
216 ax.set_xlabel('Atributo 1')
217 ax.set_ylabel('Atributo 2')
218 ax.set_zlabel('Atributo 3')
219 plt.legend()
220 plt.grid()
221 plt.show()
222
```

La gráfica obtenida es la siguiente:

Conjunto de prueba datos data\_3D clasificado con Bayes Gaussiano Naive

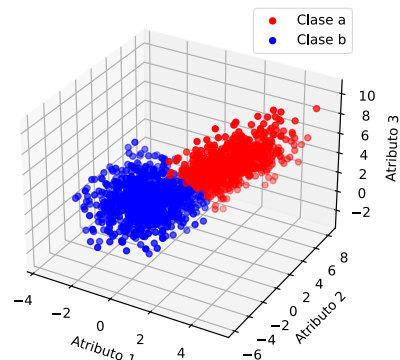


Figura 16. Conjunto de prueba clasificado con Bayes Gaussiano Naive.

- iii. Determine el error de clasificación sobre el conjunto de prueba:

Se determinó el error de clasificación, encontrando la cantidad de clasificaciones mal realizadas dividido el número de elementos clasificados.

```
223 # Cálculo error porcentual de la clasificación
224 error_bayes_naive = 100*sum(y_bayes_naive != y_real)/len(y_real)
225 print(f"Error Bayes Gaussiano Naive: {error_bayes_naive[0]:.3f}%")
```

Los errores porcentuales obtenidos en 3 ejecuciones del programa fueron:

```
Error Bayes Gaussiano Naive: 7.688%
Error Bayes Gaussiano Naive: 7.750%
Error Bayes Gaussiano Naive: 8.188%
```

#### IV. CONCLUYA SOBRE LOS RESULTADOS OBTENIDOS DE LOS CLASIFICADORES 2.C Y 2.D

Partiendo de los errores de clasificaciones encontrados, se puede decir que el clasificador Bayesiano Gaussiano funciona correctamente teniendo errores menores al 5%.

En el caso del clasificador a Bayesiano Gaussiano Naive, se siguen obteniendo buenos resultandos con errores menores al 10%. Este incremento en el error se debe a la aproximación que se hace al decir que los atributos son independientes entre sí, ya que como se pudo observar en la matriz de covarianza de la clase a, estos sí presentan una dependencia.

En general, se puede decir que ambos clasificadores tienen un buen rendimiento y al observar los resultados de clasificación.