

Actividad 2 - Componentes principales - Indicadores económicos y sociales de 96 países

- Frida Cano Falcón - A01752953

PARTE I

A partir de los datos sobre indicadores económicos y sociales de 96 países (datos: paises\_mundo.csv Download paises\_mundo.csv) hacer un análisis de Componentes principales a partir de la matriz de varianzas-covarianzas y otro a partir de la matriz de correlaciones, comparar los resultados y argumentar cuál es mejor según los resultados obtenidos.

Conexión al directorio y librerías

```
In [ ]: from google.colab import drive
drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

In [ ]: %cd "/content/drive/MyDrive/7mo Semestre/EstadísticaAvanzada"
ls
/content/drive/MyDrive/Semestres/7mo Semestre/EstadísticaAvanzada
Act1_NormalVariada_A01752953
Act2_ComponentesPrincipales_A01752953.ipynb

In [ ]: import numpy as np
from scipy.stats import multivariate_normal
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
import pandas as pd
import seaborn as sb
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler

In [ ]: # Cargar los datos a un dataframe
df = pd.read_csv('paises_mundo.csv')
df.head()

Out [ ]:
CrecPobl  MortInf  PorcMujeres  PNB95  ProdElec  LinTelf  ConsAgua  PropBosq  PropDefor  ConsEner  EmisCO2
0      1.0      30          41      2199      3903      12          94          53          0.0      341          1.2
1      3.0     124          46     4422      955          6          57          19          0.7          89          0.5
2      4.3      21          13     13540     91019          96          497          1          0.0      4566         13.1
3      2.5      34          24      44609     19883          42          180          2          0.8          906          3.0
4      1.3      22          31     278431     65962         160         1043          22          0.1         1504          3.5
```

```
In [ ]: covariance_matrix = df.cov()
covariance_matrix_df = pd.DataFrame(covariance_matrix)
covariance_matrix_df

Out [ ]:
CrecPobl  MortInf  PorcMujeres  PNB95  ProdElec  LinTelf  ConsAgua  PropBosq  PropDefor  ConsEner  EmisCO2
0      1.0      30          41      2199      3903      12          94          53          0.0      341          1.2
1      3.0     124          46     4422      955          6          57          19          0.7          89          0.5
2      4.3      21          13     13540     91019          96          497          1          0.0      4566         13.1
3      2.5      34          24      44609     19883          42          180          2          0.8          906          3.0
4      1.3      22          31     278431     65962         160         1043          22          0.1         1504          3.5
```

```
In [ ]: # Calcular la matriz de correlaciones
correlation_matrix = df.corr()
sb.heatmap(correlation_matrix, cmap = 'Blues', annot = True)

Out [ ]: <Axes: >

CrecPobl  MortInf  PorcMujeres  PNB95  ProdElec  LinTelf  ConsAgua  PropBosq  PropDefor  ConsEner  EmisCO2
CrecPobl  1.000000  0.550000  0.560000  0.320000  -0.300000  0.669016  0.200000  -0.300000  -0.180000
MortInf   0.550000  1.000000  0.033032  0.240000  -0.170000  0.070025  0.250000  -0.620000  -0.570000
PorcMujeres -0.560000  0.033032  1.000000  0.140000  0.190000  0.200000  0.310000  0.023000  0.150000
PNB95     -0.320000  0.240000  0.140000  1.000000  0.470000  0.085005  0.190000  0.280000  0.210000
ProdElec  -0.300000  -0.170000  0.190000  0.470000  1.000000  0.180000  0.260000  0.170000  0.230000
LinTelf    -0.560000  0.250000  0.200000  0.280000  0.170000  1.000000  0.063000  0.780000  0.620000
ConsAgua   -0.669016  0.310000  0.085005  0.190000  0.260000  0.063000  1.000000  0.110000  0.160000
PropBosq   -0.160000  0.200000  0.055000  0.020000  0.063000  0.110000  0.110000  1.000000  0.120000
PropDefor  -0.200000  0.230000  0.190000  0.170000  0.380000  0.087009  0.110000  0.350000  1.000000
ConsEner   -0.300000  -0.620000  0.150000  0.280000  0.780000  0.110000  0.110000  0.350000  0.860000
EmisCO2    -0.180000  -0.570000  0.047000  0.210000  0.200000  0.120000  0.120000  0.370000  0.880000
```

```
In [ ]: # Calcular los valores y vectores propios de cada matriz.
# Calcular los autovalores y autovectores Covarianza
eigenvalues_cov, eigenvectors_cov = np.linalg.eig(covariance_matrix)
# Calcular los autovalores y autovectores Correlación
eigenvalues_corr, eigenvectors_corr = np.linalg.eig(correlation_matrix)

[[-1.65816773e-06  4.70678468e-07 -1.26373574e-04 -1.92840781e-05  5.53739797e-03  1.24345612e-02 -5.35989839e-03 -8.39981843e-02  6.77857790e-02  9.07289730e-01 -1.15809996e-01]
 [-0.84013855e-05 -1.77425442e-05 -8.22582128e-03 -2.49325727e-02  9.44038283e-02  9.91751599e-01 -2.25881958e-02 -7.89112785e-02  1.63735920e-02  2.00240992e-02 -4.26487193e-04]
 [ 5.73909610e-06 -1.08454281e-05 -1.31814809e-04  5.53830717e-03  3.14036410e-02  8.55292544e-02  1.13648880e-01  9.85649883e-01  1.48846382e-02  8.34524274e-02  3.24146559e-03]
 [ 8.88037597e-01  4.59763191e-01 -2.60228711e-03 -3.89356789e-04  3.32740919e-04  8.62180481e-06  7.56647936e-06  1.21724819e-05  3.97148710e-07  2.73390791e-07  4.27445056e-07]
 [ 4.59763574e-01  8.88048472e-01 -5.69489558e-04  1.09630588e-03 -2.07619191e-04  1.85548896e-05 -1.54485759e-05 -2.58499898e-05 -1.85947801e-06 -2.08685783e-07 -1.35388952e-06]
 [ 3.58434128e-04  4.01617912e-04  6.19424889e-02 -7.64117441e-03 -9.92140409e-01  8.10902152e-02 -4.74688183e-02 -3.41681151e-02  5.37954917e-03  4.94439594e-04 -3.40942294e-03]
 [ 2.62598015e-04 -1.12211782e-03  4.01453227e-02 -9.99411424e-01 -5.77951444e-03 -1.08722945e-03  8.86329425e-03 -4.69073884e-03 -7.86526863e-05  4.78841039e-04  3.62142471e-05]
 [ 4.88956383e-06  7.79884284e-06 -1.27193178e-03  6.43579663e-03 -4.19316151e-02  1.72184892e-02  9.92053003e-01 -1.18963839e-01 -1.41656633e-03 -3.74897888e-03  5.89175814e-03]
 [-1.73702529e-06  2.35888615e-07 -1.91617666e-04  4.04379633e-05  1.89007507e-03  1.75866759e-03 -7.45542659e-03  1.81144537e-02 -1.28383989e-01 -1.05293378e-01 -9.85931887e-01]
 [ 2.64715880e-04 -1.12078218e-04  4.97231506e-01  3.97356775e-02  6.52739475e-02  2.63867318e-03  3.78478929e-03  1.26120513e-03 -2.8293981e-03  5.98624159e-05  2.67261847e-04]
 [ 4.64372357e-06 -1.51575121e-06  2.86799405e-03 -5.62684659e-05  4.29871157e-03 -1.87798798e-02 -1.70913698e-03 -5.28482298e-03  9.89152936e-01 -8.22137894e-02 -1.20851858e-01]]
```

```
In [ ]: # Proporción de varianza explicada por cada componente
var_total = np.trace(covariance_matrix) # Varianza Total sum(diag(covariance_matrix))
prop_var = eigenvalues_cov / var_total
print("Proporción de varianza: ", prop_var)

Proporción de varianza: [ 0.83454311e-01  9.64729842e-02  6.79580362e-05  4.55456679e-06  1.78242937e-07  7.53891641e-09  5.31773802e-09  6.85776295e-10  8.50289159e-11  6.88905312e-12  2.10784322e-11]

In [ ]: # Calcular la suma acumulativa
suma_acumulativa = np.cumsum(prop_var)
print(suma_acumulativa)

[0.83454311 0.99992713 0.99999525 0.99999981 0.99999999 0.99999999
 1.         1.         1.         1.         1.]
```

Para determinar qué componentes son los más importantes, observamos las proporciones de varianza explicada acumulada. Los componentes con una mayor proporción acumulada explican una mayor cantidad de variabilidad en los datos:

CrecPobl	MortInf	PorcMujeres	PNB95	ProdElec	LinTelf	ConsAgua	PropBosq	PropDefor	ConsEner	EmisCO2
0.903	0.09	0.00006	0.000000455	0.00000178	0.0000000075	0.0000000063	0.0000000006	0.00000000085	0.000000000069	0.000000000021

En este caso se reconoce que los componentes con mayor proporción acumulada son CrecPobl y MortInf

PARTE II

Obtenga las gráficas de respectivas con S (matriz de varianzas-covarianzas) y con R (matriz de correlaciones) de las dos primeras componentes e interprete los resultados en término de agrupación de variables (puede ayudar "Índice de riqueza", "Índice de ruralidad")

```
In [ ]: # Estándarizar los datos (importante antes de realizar PCA)
scaler = StandardScaler()
datos_escalados = scaler.fit_transform(df[['CrecPobl', 'MortInf']])

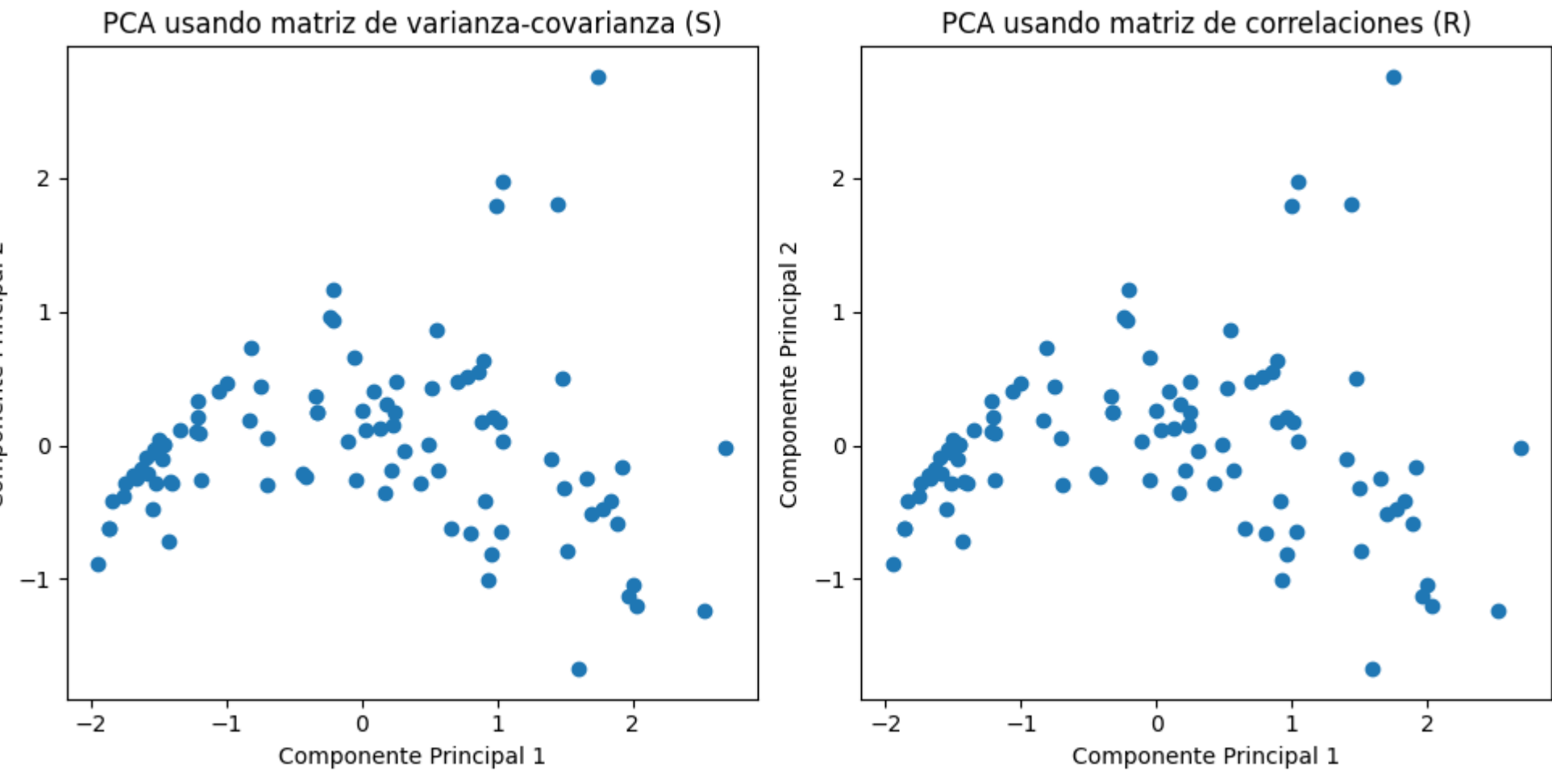
In [ ]: # Realizar PCA utilizando la matriz de varianza-covarianza (S)
pca_cov = PCA(n_components=2)
componentes_principales_cov = pca_cov.fit_transform(datos_escalados)

# Realizar PCA utilizando la matriz de correlaciones (R)
pca_cor = PCA(n_components=2)
componentes_principales_cor = pca_cor.fit_transform(datos_escalados)

# Graficar las dos primeras componentes principales usando S
plt.figure(figsize=(10, 5))
plt.subplot(1, 2, 1)
plt.scatter(componentes_principales_cov[:, 0], componentes_principales_cov[:, 1])
plt.title("PCA usando matriz de varianza-covarianza (S)")
plt.xlabel("Componente Principal 1")
plt.ylabel("Componente Principal 2")

# Graficar las dos primeras componentes principales usando R
plt.subplot(1, 2, 2)
plt.scatter(componentes_principales_cor[:, 0], componentes_principales_cor[:, 1])
plt.title("PCA usando matriz de correlaciones (R)")
plt.xlabel("Componente Principal 1")
plt.ylabel("Componente Principal 2")

plt.tight_layout()
plt.show()
```



PARTE III

Explore los siguientes gráficos relativos al problema y Componentes Principales y dé una interpretación de cada gráfico.

```
In [ ]: scaler = StandardScaler()
datos_escalados = scaler.fit_transform(df)

# Realizar PCA
pca = PCA()
componentes_principales = pca.fit_transform(datos_escalados)

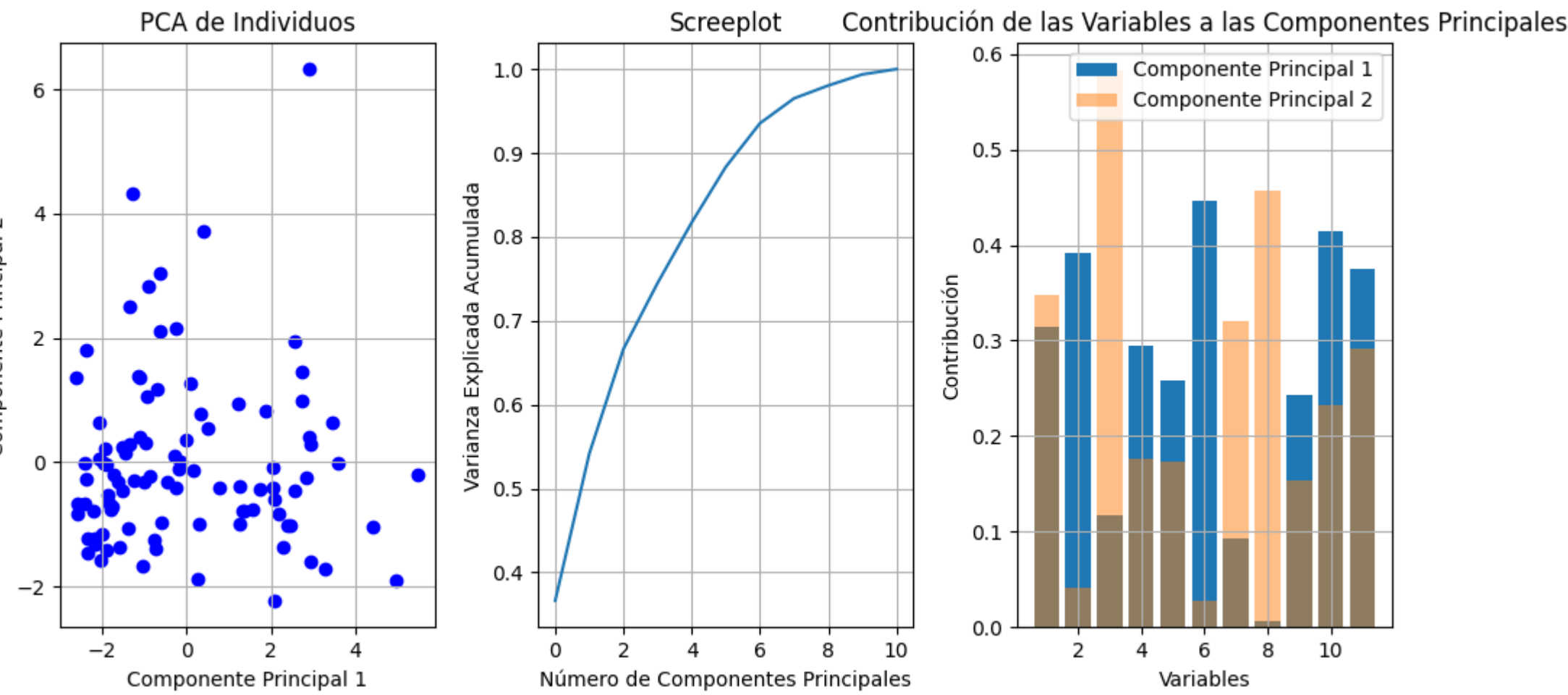
# Visualizar individuos (puedes usar la biblioteca matplotlib)
plt.figure(figsize=(10, 5))
plt.subplot(1, 3, 1)
plt.scatter(componentes_principales[:, 0], componentes_principales[:, 1], c="blue")
plt.xlabel("Componente Principal 1")
plt.ylabel("Componente Principal 2")
plt.title("PCA de Individuos")
plt.grid()

# Visualizar el screeplot (varianza explicada acumulada)
plt.subplot(1, 3, 2)
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.xlabel("Número de Componentes Principales")
plt.ylabel("Varianza Explicada Acumulada")
plt.title("Screeplot")
plt.grid()

# Calcular la contribución de las variables a las componentes principales
contribuciones = pca.components_

# Visualizar las contribuciones (reemplaza esto con tu propio código para graficar)
# Por ejemplo, puedes crear un gráfico de barras para mostrar las contribuciones de las variables a las componentes principales
variables = range(1, df.shape[1] + 1)
plt.subplot(1, 3, 3)
plt.bar(variables, np.abs(contribuciones[0]), label="Componente Principal 1")
plt.bar(variables, np.abs(contribuciones[1]), label="Componente Principal 2", alpha=0.5)
plt.xlabel("Variables")
plt.ylabel("Contribución")
plt.title("Contribución de las Variables a las Componentes Principales")
plt.legend()
plt.grid()

plt.tight_layout()
plt.show()
```



De acuerdo a los gráficos obtenidos:

- Gráfico de dispersión de individuos (PCA de Individuos):
  - Este gráfico muestra la proyección de los individuos (muestras o observaciones) en el espacio de las dos primeras componentes principales.
  - Cada punto en el gráfico representa un individuo.
  - El eje x representa el valor de la primera componente principal, mientras que el eje y representa el valor de la segunda componente principal.
  - Puedes observar cómo los individuos se agrupan o dispersan en función de su posición en el espacio de las componentes principales.
  - Esto puede ayudar a identificar patrones de similitud o diferencia entre las observaciones.
- Screeplot (Varianza Explicada Acumulada):
  - Este gráfico muestra la varianza explicada acumulada en función del número de componentes principales.
  - Cada punto en el gráfico representa la proporción acumulada de varianza explicada a medida que se agregan más componentes principales.
  - El eje x muestra el número de componentes principales, y el eje y muestra la varianza explicada acumulada.
  - Puedes utilizar este gráfico para determinar cuántas componentes principales retener. Generalmente, se buscan "codos" en el gráfico para decidir el número apropiado de componentes principales a conservar. En este caso, parece que las dos primeras componentes explican la mayoría de la varianza.
- Gráfico de Contribución de Variables a las Componentes Principales:
  - Este gráfico muestra la contribución de las variables originales a las dos primeras componentes principales.
  - Cada barra en el gráfico representa una variable, y la altura de la barra indica la magnitud de esa variable a las componentes principales.
  - Puedes observar qué variables tienen las mayores contribuciones a cada componente principal.
  - Las variables con contribuciones altas (positivas o negativas) son las que más influyen en la dirección y magnitud de las componentes principales. Esto puede ayudarte a identificar qué variables están más relacionadas con las dimensiones subyacentes de tus datos.

