

In [1]:

```
INFORME

CLEANING BANK MARKETING CAMPAIGN DATA

Noviembre 15, 2024

1.1.1. Importar paquetes

In [14]:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.compose import ColumnTransformer
from sklearn.cluster import DBSCAN
from mpl_toolkits.mplot3d import Axes3D
from sklearn.decomposition import PCA

1.2.2. Cargar los datos

In [15]:
df_economico = pd.read_csv('economico.csv')
df_cliente = pd.read_csv('cliente.csv')
df_campaign = pd.read_csv('campaign.csv')

In [1]:
```

```
1.3.3. Mostramos informacion de cada Data Frame

In [16]:
df_economico.info(), df_cliente.info(), df_campaign.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4188 entries, 0 to 4187
Data columns (total 3 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   client_id              4188 non-null    int64
 1   cont_precio_ide       4188 non-null    float64
 2   wallet_share_masha    4188 non-null    float64
dtypes: float64(2), int64(1)
memory usage: 260.0 KB

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4188 entries, 0 to 4187
Data columns (total 7 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   client_id              4188 non-null    int64
 1   age                   4188 non-null    int64
 2   job                   4188 non-null    object
 3   marital               4188 non-null    object
 4   education             3947 non-null    object
 5   credit_default        4188 non-null    bool
 6   mortgage              4187 non-null    bool
dtypes: bool(2), int64(1), object(4)
memory usage: 1.4+ MB

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4188 entries, 0 to 4187
Data columns (total 7 columns):
 #   Column                Non-Null Count  Dtype
---  --
 0   client_id              4188 non-null    int64
 1   number_contacts       4188 non-null    int64
 2   contact_duration      4188 non-null    object
 3   previous_campaign_contacts  4188 non-null    int64
 4   previous_outcome      4188 non-null    bool
 5   campaign_outcome      4188 non-null    bool
 6   last_contact_date     4188 non-null    object
dtypes: bool(2), int64(1), object(4)
memory usage: 1.4+ MB

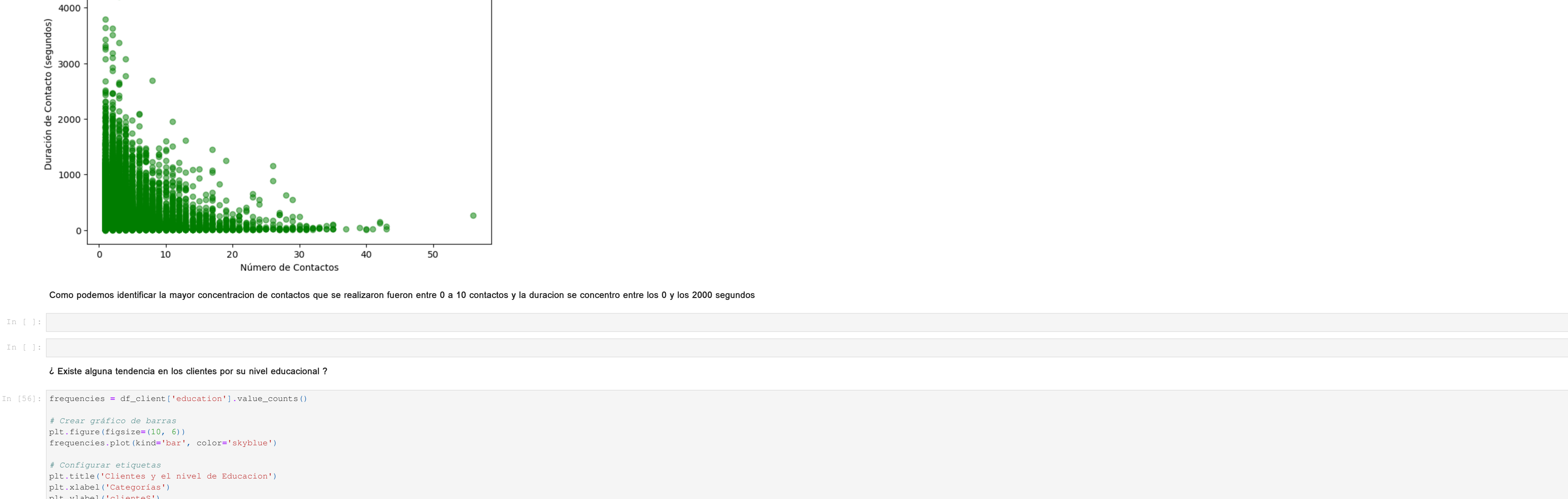
In [14]:
df_campaign.info()
Out0, Out1, Out2
```

2.1 Creamos un analisis Grafico

¿ La Matriz de datos nos ampa informacion relevante ?

```
In [15]:
df_campaign = np.random.randint(0, 10) # Matriz de 10x10 con valores aleatorios entre 0 y 1

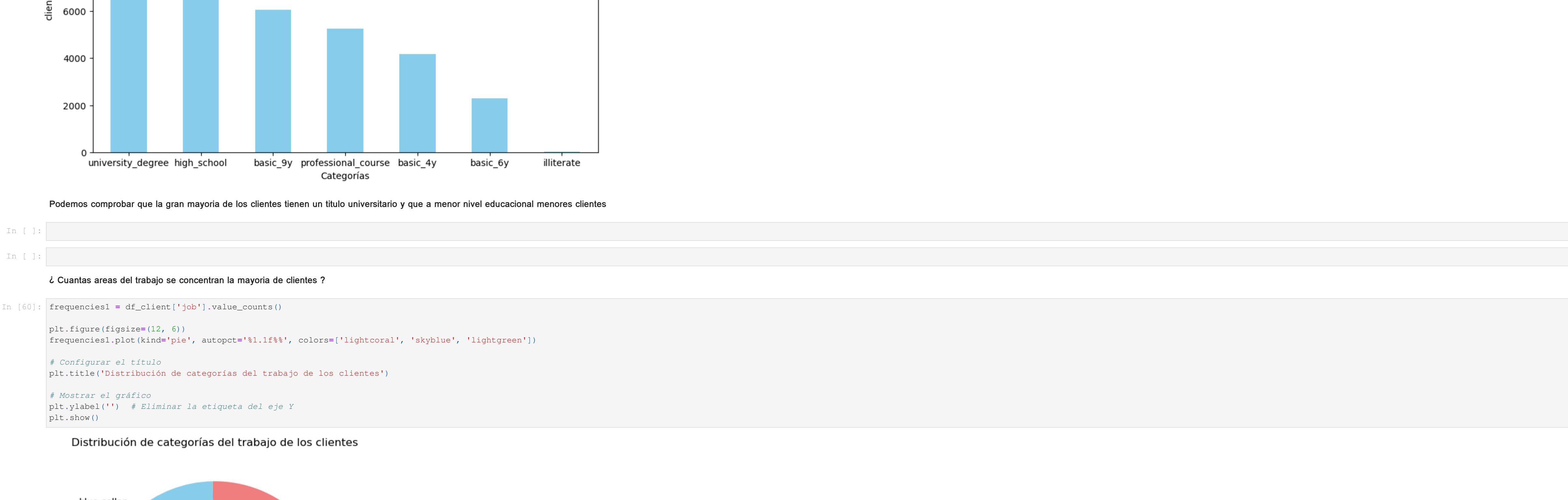
# Crear la grafica
plt.imshow(df_campaign, cmap='viridis', interpolation='none')
plt.colorbar() # Agregar una barra de colores
plt.title('Matriz de Datos')
plt.xlabel('Columnas')
plt.ylabel('Filas')
plt.show()
```



No loge identificar alguna informacion relevante asociada a la matriz generada del dataframe de la campaña

¿ Que cantidad de contactos hubieron en relacion a la duracion (segundos) ?

```
In [14]:
plt.figure(figsize=(8, 6))
plt.scatter(df_campaign['number_contacts'], df_campaign['contact_duration'], color='green', alpha=0.1)
plt.title('Duración de Contactos vs Numero de Contactos')
plt.xlabel('Numero de Contactos')
plt.ylabel('Duración de Contacto (segundos)')
plt.show()
```



Como podemos identificar la mayor concentracion de contactos que se realizaron fuera entre 0 a 10 contactos y la duracion se concentro entre los 0 y los 2000 segundos

In [1]:

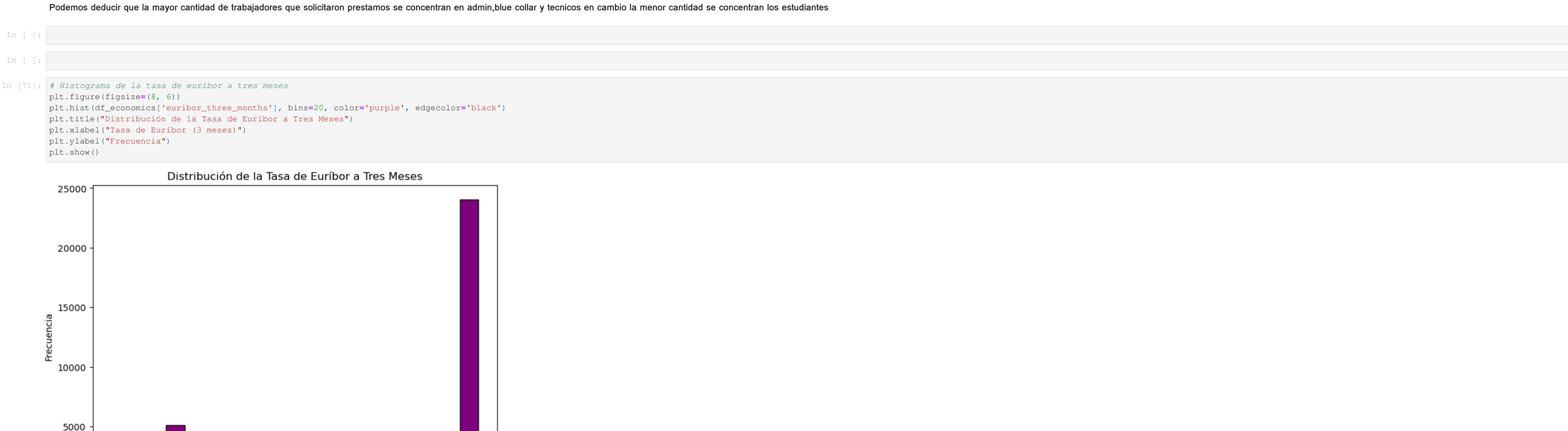
¿ Existe alguna tendencia en los clientes por su nivel educacional ?

```
In [16]:
frecuencias = df_cliente['education'].value_counts()

# Crear grafico de barras
plt.figure(figsize=(10, 6))
frecuencias.plot(kind='bar', color='skyblue')

# Configurar etiquetas
plt.title('Clientes y el nivel de Educacion')
plt.xlabel('Categorías')
plt.ylabel('Frecuencia')
plt.xticks(rotation=90) # Para que las categorias no se roten

# Mostrar el gráfico
plt.show()
```



Podemos comprobar que la gran mayoría de los clientes tienen un tibu universitario y que a menor nivel educacional menores clientes

In [1]:

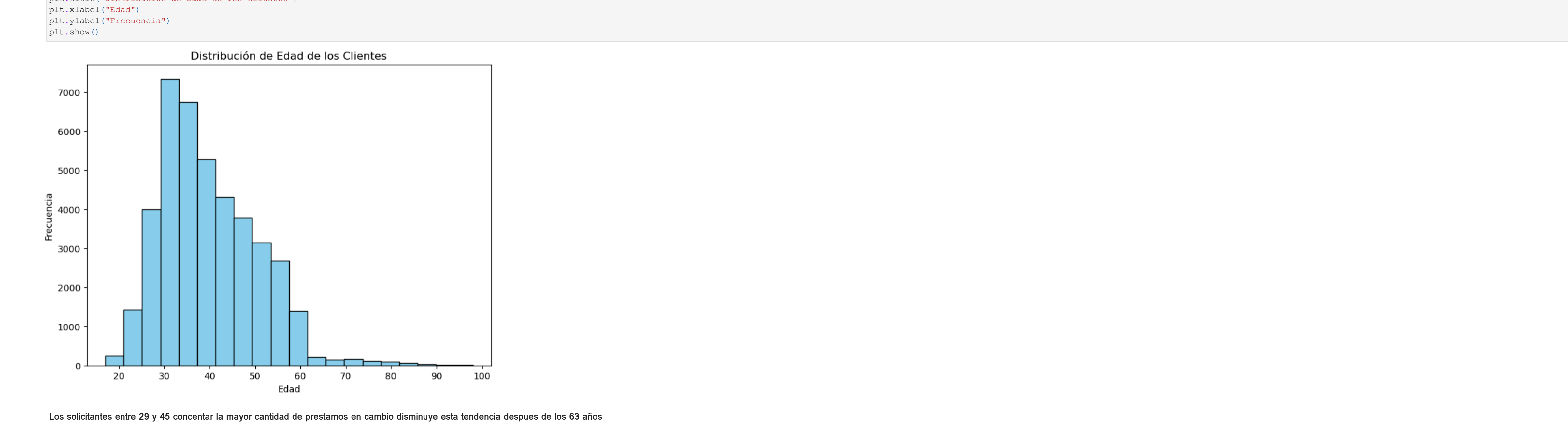
¿ Cuántas areas del trabajo se concentran la mayoría de clientes ?

```
In [16]:
frecuencias1 = df_cliente['job'].value_counts()

plt.figure(figsize=(12, 6))
frecuencias1.plot(kind='pie', autopct='%1.1f%%', colors=['lightcoral', 'skyblue', 'lightgreen'])

# Configurar el título
plt.title('Distribución de categorías del trabajo de los clientes')

# Mostrar el gráfico
plt.xlabel('') # Eliminar la etiqueta del eje x
plt.ylabel('') # Eliminar la etiqueta del eje y
plt.show()
```

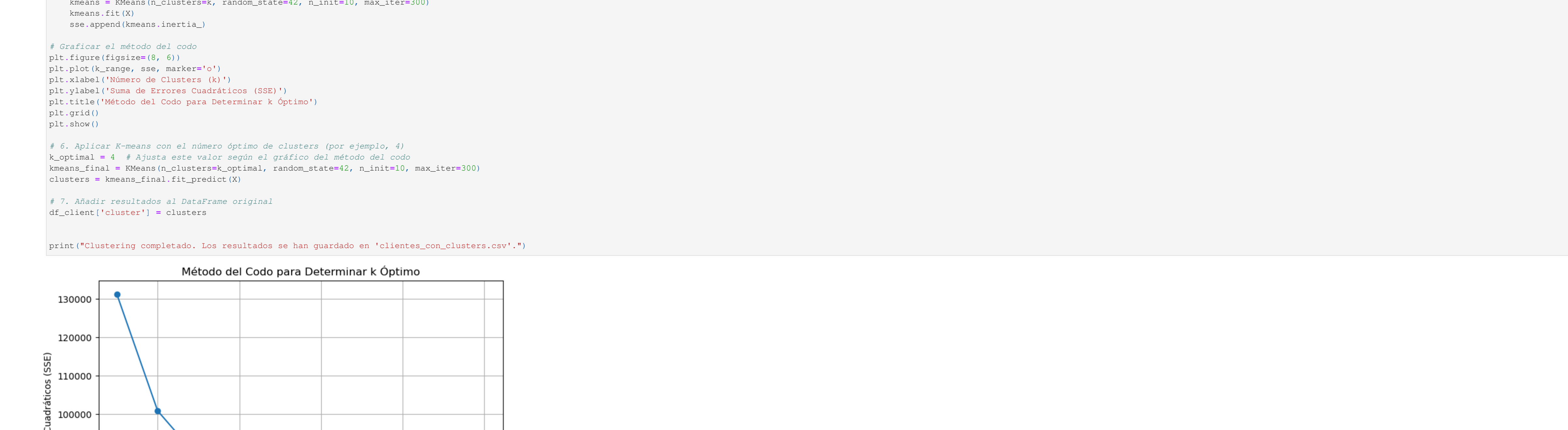


Podemos deducir que la mayor cantidad de trabajadores que solicitaron prestamos se concentran en admin, blue collar y tecnicos en cambio la menor cantidad se concentran los estudiantes

In [1]:

¿ La edad es un factor para solicitar prestamos ?

```
In [17]:
# Histograma de la tasa de euribor a tres meses
plt.figure(figsize=(8, 6))
plt.hist(df_economico['euribor_tres_meses'], bins=20, color='purple', edgecolor='black')
plt.title('Distribución de la Tasa de Euribor a Tres Meses')
plt.xlabel('Tasa de Euribor (3 meses)')
plt.ylabel('Frecuencia')
plt.show()
```



Forma de la distribución: Observa la forma de la distribución, si es simétrica, sesgada a la derecha o a la izquierda. Por ejemplo, si la mayoría de los valores están en el lado izquierdo, indica una asimetría negativa (peseo hacia valores bajos), si están en el lado derecho, muestra una asimetría positiva. Rango de tasas: Identifica el rango en el eje x para ver cómo varía la tasa de Euribor. Esto puede indicar en qué intervalos de valores es más probable que se encuentren las tasas. Frecuencia de valores específicos: Observa los picos para identificar las tasas de Euribor más comunes durante el periodo de observación

In [1]:

¿ La edad es un factor para solicitar prestamos ?

```
In [17]:
# 2. Eliminar columnas irrelevantes
df_cliente = df_cliente.drop(columns='client_id')

# 3. Preparar variables para codificación y escalado
categoricas_columnas = ['sex', 'marital', 'education']
numericas_columnas = ['age']

# 4. Codificar variables categoricas y escalar numericos
preprocesador = ColumnTransformer(
    transformers=[
        ('cat', StandardScaler(), numericas_columnas),
        ('out', OneHotEncoder(), categoricas_columnas)
    ]
)

# Transformar los datos
X = preprocesador.fit_transform(df_cliente)

# 5. Método del codo
nro = 1

# Rango = range(1, 11) # Probar entre 1 y 10 clusters
for k in k_range:
    kmeans = KMeans(n_clusters=k, random_state=42, n_init=10, max_iter=300)
    kmeans.fit(X)
    ss.append(kmeans.inertia_)

# Graficar el método del codo
plt.figure(figsize=(8, 6))
plt.plot(k_range, ss, marker='x')
plt.xlabel('Numero de Clusters (k)')
plt.ylabel('Suma de Errores Cuadráticos (SSE)')
plt.title('Método del Codo para Determinar k Óptimo')
plt.grid()

# 6. Aplicar Kmeans con el número óptimo de clusters (por ejemplo, 4)
k_optimal = 4 # Ajusta este valor según el gráfico del método del codo
kmeans_final = KMeans(n_clusters=k_optimal, random_state=42, n_init=10, max_iter=300)
clusters = kmeans_final.fit_predict(X)

# 7. Guardar resultados al DataFrame original
df_cliente['cluster'] = clusters

print('Clustering completado. Los resultados se han guardado en 'clientes_con_clusters.csv'.')
```

Clustering completado. Los resultados se han guardado en 'clientes\_con\_clusters.csv'. Según el resultado del metodo del codo vamos a usar 3 cluster

In [1]:

```
In [18]:
# 3. Crear 3 clusters
k_optimal = 3 # Cambiado para usar 3 clusters
kmeans_final = KMeans(n_clusters=k_optimal, random_state=42, n_init=10, max_iter=300)
clusters = kmeans_final.fit_predict(X)

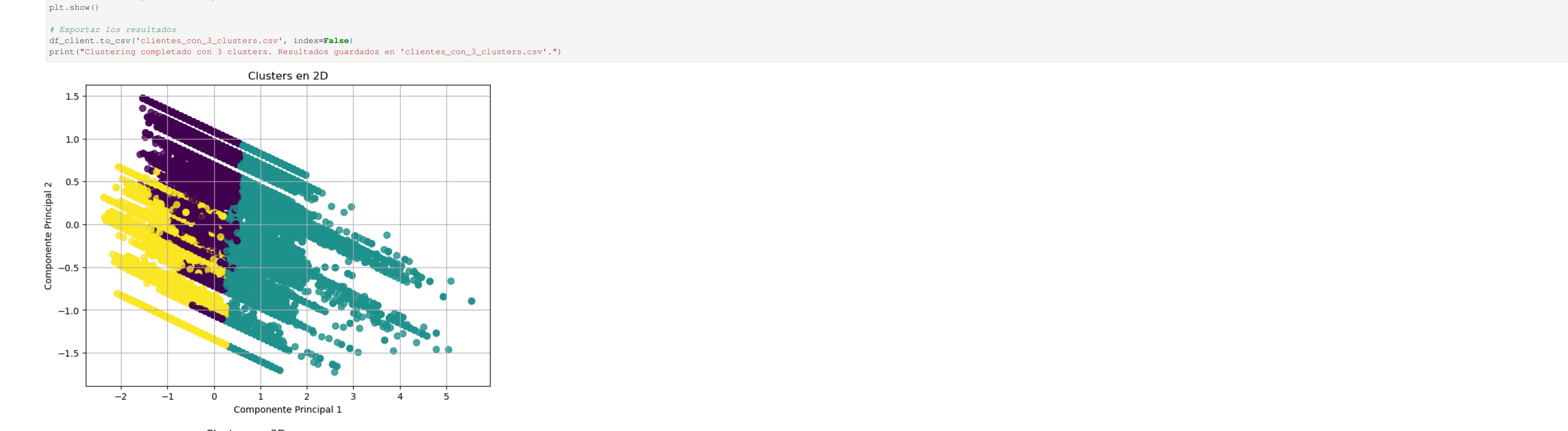
# 4. Guardar resultados al DataFrame original
df_cliente['cluster'] = clusters

# 5. Reducir la dimensionalidad para graficar (PCA a 2D y 3D)
pca = PCA(n_components=2)
X_2d = pca.fit_transform(X)
X_3d = pca.fit_transform(X)

# Graficar en 2D
plt.figure(figsize=(8, 6))
plt.scatter(X_2d[:, 0], X_2d[:, 1], c=clusters, cmap='viridis', s=50, alpha=0.8)
plt.title('Clusters en 2D')
plt.xlabel('Componente Principal 1')
plt.ylabel('Componente Principal 2')
plt.grid()
plt.show()

# Graficar en 3D
from mpl_toolkits.mplot3d import Axes3D
fig = plt.figure(figsize=(10, 8))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(X_3d[:, 0], X_3d[:, 1], X_3d[:, 2], c=clusters, cmap='viridis', s=50, alpha=0.8)
ax.set_xlabel('Componente Principal 1')
ax.set_ylabel('Componente Principal 2')
ax.set_zlabel('Componente Principal 3')
plt.show()

# Mostrar los resultados
df_cliente_con_clusters_3d = df_cliente
print('Clustering completado con 3 clusters. Resultados guardados en 'clientes_con_clusters.csv'.')
```



Clustering completado con 3 clusters. Resultados guardados en 'clientes\_con\_clusters.csv'.

In [18]:

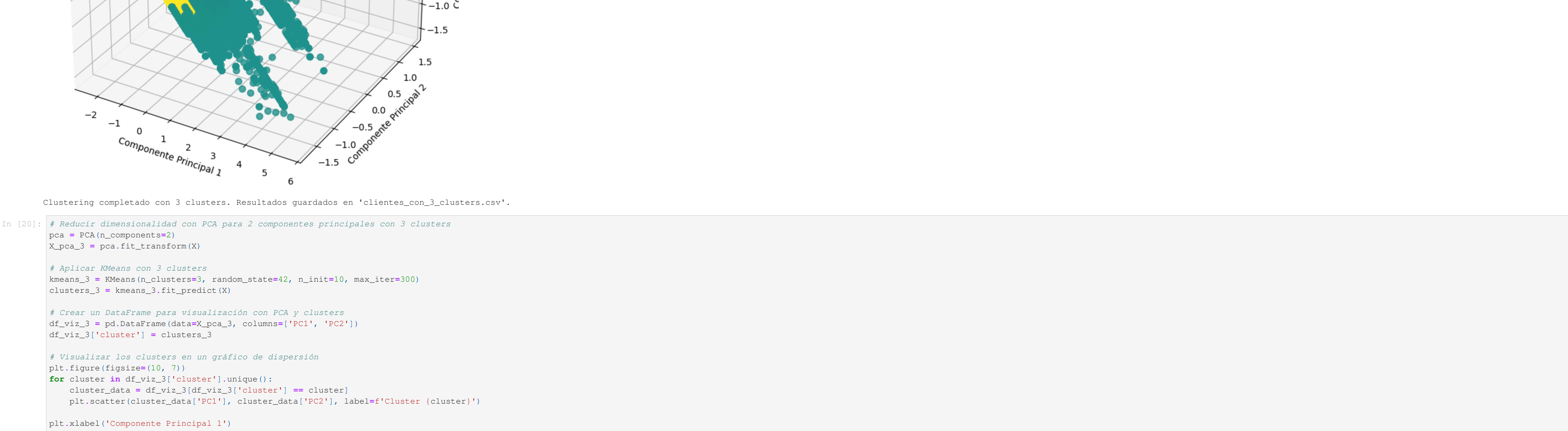
```
# 5. Reducir dimensionalidad con PCA para 2 componentes principales con 3 clusters
pca = PCA(n_components=2)
X_2d = pca.fit_transform(X)

# Aplicar DBSCAN con 3 clusters
dbSCAN = DBSCAN(eps=0.5, min_samples=5, n_init=10, max_iter=300)
clusters_3 = dbSCAN.fit_predict(X)

# Crear un DataFrame para visualización con PCA y clusters
df_cliente_3d = df_cliente.copy()
df_cliente_3d['cluster'] = clusters_3

# Visualizar los clusters en un gráfico de dispersión
plt.figure(figsize=(10, 8))
for cluster in df_cliente_3d['cluster'].unique():
    cluster_data = df_cliente_3d[df_cliente_3d['cluster'] == cluster]
    plt.scatter(cluster_data['PC1'], cluster_data['PC2'], label=f'Cluster {cluster}')

plt.xlabel('Componente Principal 1')
plt.ylabel('Componente Principal 2')
plt.grid()
plt.show()
```



Clustering completado con 3 clusters. Resultados guardados en 'clientes\_con\_clusters.csv'.

In [18]:

```
# 5. Reducir dimensionalidad con PCA para 2 componentes principales con 3 clusters
pca = PCA(n_components=2)
X_2d = pca.fit_transform(X)

# Aplicar DBSCAN con 3 clusters
dbSCAN = DBSCAN(eps=0.5, min_samples=5, n_init=10, max_iter=300)
clusters_3 = dbSCAN.fit_predict(X)

# Crear un DataFrame para visualización con PCA y clusters
df_cliente_3d = df_cliente.copy()
df_cliente_3d['cluster'] = clusters_3

# Visualizar los clusters en un gráfico de dispersión
plt.figure(figsize=(10, 8))
for cluster in df_cliente_3d['cluster'].unique():
    cluster_data = df_cliente_3d[df_cliente_3d['cluster'] == cluster]
    plt.scatter(cluster_data['PC1'], cluster_data['PC2'], label=f'Cluster {cluster}')

plt.xlabel('Componente Principal 1')
plt.ylabel('Componente Principal 2')
plt.grid()
plt.show()
```

