	PCA (Principal Component Analysis)  PCA es el método más utilizado para reducción de dimensionalidad. Se utiliza para reducir el número de atributos o variables de un conjunto de datos para facilitar el análisis. PCA considera como "información" la varianza de cada atributo y busca reducir duplicidad de información en términos de variables correlacionadas.
Γ	#importamos librerías import pandas as pd import numpy as np import seaborn as sns import matplotlib.pyplot as plt  NOTA: La base de datos correspondiente a la actividad anterior de limpieza de datos se guardó en el archivo 'clients_clean.csv'
[]:	#Creamos dataframe del archivo csv df = pd.read_csv('clients_clean.csv', index_col=0)    ID   X1   X2   X3   X4   X5   X6   X7   X8   X9     X15   X16   X17   X18   X19   X20   X21   X22   X23   Y
	2 3 9000 2.0 2.0 2.0 4.0 4.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1.0 1
	La descripción de la base de datos la encontramos aqui: https://raw.githubusercontent.com/PosgradoMNA/Actividades_Aprendizaje-/main/dataset_info.txt  Paso 1: Determine el número mínimo de componentes principales que representan la mayor parte de la variación en sus datos  • Utilice la proporción acumulada de la varianza que explican los componentes para determinar la cantidad de varianza que explican los componentes principales.
[]:	Recordamos que PCA funciona solo con variables numericas, por lo que descartamos las columnas con datos categóricos.  # Creamos una copia del dataframe, eliminando variables categoricas ndf = df.copy() ndf.drop(['ID', 'X2', 'X3', 'X4', 'X6', 'X7', 'X8', 'X9', 'X10', 'X11', 'Y'],
	0         20000         24.0         3913.0         3102.0         689.0         0.0         0.0         0.0         689.0         0.0
	2995         220000         39.0         188948.0         192815.0         208365.0         88004.0         31237.0         15980.0         8500.0         20000.0         5003.0         3047.0         5000.0         1000.0           29996         150000         43.0         1683.0         1828.0         3502.0         8979.0         5190.0         0.0         1837.0         3526.0         8998.0         129.0         0.0         0.0           29997         30000         37.0         3565.0         3356.0         2758.0         20878.0         20582.0         19357.0         0.0         0.0         22000.0         4200.0         2000.0         3100.0           29998         80000         41.0         1645.0         78379.0         76304.0         52774.0         11855.0         48944.0         85900.0         3409.0         1178.0         1926.0         52964.0         1804.0           29999         50000         46.0         47929.0         48905.0         49764.0         36535.0         32428.0         15313.0         2078.0         1800.0         1430.0         1000.0         1000.0         1000.0
	Ahora que tenemos solo las columnas para nuestro análisis, procedemos a revisar la correlación entre variables.  Para visualizar mejor los datos, creamos un mapa de calor en donde los valores en color mas claro (cercanos a 1) indican mayor correlación entre pares de variables y los colores mas oscuro (cercanos a 0) indican poca o ninguna correlación.  # Creamos matríz de correlación con mapa de calor
[]:	<pre>sns.set(rc={'figure.figsize':(12,10)}) corr_mat = ndf.corr(method='pearson') sns.heatmap(corr_mat, annot=True)  <axessubplot:>  1  0.14  0.29  0.28  0.28  0.29  0.3  0.29  0.2  0.18  0.21  0.2  0.22  0.22  Ø  0.14  1  0.056  0.054  0.052  0.049  0.048  0.026  0.022  0.029  0.021  0.023  0.019</axessubplot:></pre>
	0.29 0.056 1 0.95 0.89 0.86 0.83 0.8 0.15 0.1 0.16 0.16 0.17 0.18  0.28 0.28 0.054 0.95 1 0.93 0.89 0.86 0.83 0.85 0.25 0.32 0.13 0.14 0.18  0.28 0.29 0.052 0.86 0.89 0.93 1 0.94 0.9 0.23 0.21 0.3 0.14 0.16 0.18  -0.8  -0.8  -0.8  -0.8  -0.8  -0.8  -0.8  -0.8  -0.8  -0.8
	\tilde{\text{X}}       0.3       0.049       0.83       0.86       0.88       0.94       1       0.95       0.22       0.18       0.25       0.29       0.14       0.16         \tilde{\text{X}}       0.29       0.048       0.8       0.83       0.85       0.9       0.95       1       0.2       0.17       0.23       0.25       0.31       0.13         \tilde{\text{X}}       0.18       0.022       0.11       0.12       0.22       0.22       0.22       0.12       0.18       0.17       0.29       1       0.24       0.18       0.16
	0.21 0.029 0.16 0.15 0.13 0.3 0.25 0.23 0.25 0.24 1 0.22 0.16 0.16  0.2 0.021 0.16 0.15 0.14 0.14 0.29 0.25 0.2 0.18 0.22 1 0.15 0.16  0.2 0.023 0.17 0.16 0.18 0.16 0.14 0.31 0.15 0.18 0.16 0.15 1 0.15  0.22 0.019 0.18 0.17 0.18 0.18 0.16 0.13 0.19 0.16 0.16 0.15 1  X1 X5 X12 X13 X14 X15 X16 X17 X18 X19 X20 X21 X22 X23
	De las variables X12 a X17 tienenen mayor correlación entre sí, lo cual tiene lógica ya que el monto de deuda es similar de un mes a otro.  • X12 a X17 = monto de deuda de Septiembre a Abril  La variable X5 es la que menos correlación tiene con el resto de las variables, de lo cual podemos inferir que en este caso la edad del cliente no tiene tanta influencia en su comportamiento financiero.
	# importamos modulos de scikit-learn from sklearn import preprocessing from sklearn.preprocessing import StandardScaler from sklearn.decomposition import PCA  Antes de iniciar con el PCA es necesario estandarizar los datos, de lo contrario se le asignará mayor importancia a las variables cuya escala de medición sea mayor, lo cual no es necesarian
	cierto.  Usamos StandardScaler para estandarizar los datos. A esta técnica se le conoce como Normalización de puntuación Z  Este proceso estandariza cada valor de manera que:  • La media de todos los valores = 0  • La desviación estándar = 1
	Cada puntuación z nos dice a cuántas desviaciones estándar se encuentra un valor individual de la media Formula: $x_{new} = \frac{x_{old} - \mu}{\sigma}$ # Usamos StandardScaler para estandarizar los datos scaler = StandardScaler()
	# USAMOS StandardScaler   para estandar1221 108 datos   scaler = StandardScaler()   ndf_scaled = pd.DataFrame(scaler.fit_transform(ndf), columns = ndf.columns)   columns = ndf.columns)
	1         -0.3661         -1.0289         -0.6602         -0.6406         -0.6233         -0.6079         -0.6013         -0.3419         -0.2136         -0.2400         -0.2442         -0.3141         -0.1809           2         -0.5973         -0.1610         -0.2994         -0.4950         -0.4837         -0.4513         -0.3947         -0.2502         -0.1919         -0.2400         -0.2442         -0.2487         -0.0122           3         -0.9056         0.1645         -0.0582         -0.0141         0.0319         -0.2337         -0.1881         -0.1593         -0.2211         -0.1694         -0.2286         -0.2379         -0.2442         -0.2371           4         -0.9056         2.3344         -0.5795         -0.6124         -0.1623         -0.3484         -0.3344         -0.2211         1.3349         0.2712         0.2664         -0.2690         -0.2552
[]:[	29994 -1.0597 0.1645 -0.6482 -0.6450 -0.6395 -0.3494 -0.3260 -0.3306 -0.3419 -0.2570 0.9528 -0.0400 -0.1832 -0.1190 29995 -0.6744 0.5985 -0.6743 0.4097 0.4217 0.1468 -0.4696 0.1669 4.8445 -0.1091 -0.2299 -0.1851 3.1524 -0.1919 29996 -0.9056 1.1410 -0.0454 -0.0047 0.0387 -0.1058 -0.1310 -0.3986 -0.2164 -0.1789 -0.2155 -0.2442 -0.2487 -0.2371 29997 rows × 14 columns  # Calculamos la media y desviación est. de cada variable
[]:	<pre>ndf_scaled_dict = {'mean':ndf_scaled[ndf_scaled.columns].mean().round(4),</pre>
[]:	Observamos ahora que la escala de todas las variables es la misma, la media = 0 y desviación est. = 1  Procedemmos con el PCA  # Realizamos PCA  pcs = PCA()  pc comp = pcs fit transform(ndf scaled)
	<pre>p_comp = pcs.fit_transform(ndf_scaled)  # Creamos una tabla con los datos de los componentes principales pcsSummary_df = pd.DataFrame({'Desviacion estandar': np.sqrt(pcs.explained_variance_),</pre>
[]:	pcsSummary_df_round(4)    PC1   PC2   PC3   PC4   PC5   PC6   PC7   PC8   PC9   PC10   PC11   PC12   PC13   PC14
[]:[	Como podemos observar en la <i>Proporción acumulada</i> , los primeros 7 componentes explican poco mas del 86% de la varianza asociada al total de las 14 variables.  print('La varianza acumulada de los primeros 7 componentes es:',
	<pre>PC_components = np.arange(pcs.n_components_) + 1 sns.set(style = 'whitegrid', font_scale = 1.2) fig, ax = plt.subplots(figsize=(12, 7)) sns.barplot(x = PC_components, y = pcs.explained_variance_ratio_, color = 'c') sns.lineplot(x = PC_components-1,</pre>
	plt.title('Scree Plot') plt.xlabel('N-Componente Principal') plt.ylabel('Varianza Explicada') plt.ylim(0, 1) plt.show()  Scree Plot
	0.8 Explicada o.6
	1 2 3 4 5 6 7 8 9 10 11 12 13 14  Comprobamos que el total de varianza de las variables originales (estandarizadas) es igual al total de varianza de los componentes principales.  Esto nos indica que si consideramos los 14 PCs no habría pérdida de información.  # Creamos dataframe con la información de los PC's
	pc_df = pd.DataFrame(p_comp) print('Varianza total de variables originales =', ndf_scaled.var().sum().round(4)) print('Varianza total componentes principales =', pc_df.var().sum().round(4))  Varianza total de variables originales = 14.0005 Varianza total componentes principales = 14.0005  Conclusión Paso 1:  De acuerdo a la varianza acumulada de los componentes principales y estableciendo un criterio mínimo del 85%, determinamos que:
	Los primeros 7 componentes principales cumplen con el criterio, explicando el 86.72% de la varianza.  Paso 2: Interprete cada componente principal en términos de las variables originales
	<ul> <li>Examine la magnitud y la dirección de los coeficientes de las variables originales.</li> </ul>
	<ul> <li>Examine la magnitud y la dirección de los coeficientes de las variables originales.</li> <li>Nota: Cuanto mayor sea el valor absoluto del coeficiente, más importante será la variable correspondiente en el cálculo del componente.</li> <li>Es importante mencionar que cada componente principal representa un vector con la combinación de las contribuciones de varianza de cada variable original.</li> </ul>
[]:	<ul> <li>Examine la magnitud y la dirección de los coeficientes de las variables originales.</li> <li>Nota: Cuanto mayor sea el valor absoluto del coeficiente, más importante será la variable correspondiente en el cálculo del componente.</li> <li>Es importante mencionar que cada componente principal representa un vector con la combinación de las contribuciones de varianza de cada variable original.</li> <li>El PC1 captura la mayor cantidad de varianza (como se observa en la scree plot), sucesivamente el PC2 captura la mayor varianza posible bajo la condición de que este es ortogonal al es decir que no esta correlacionado con PC1. Los demás PCs siguen este mismo esquema.</li> <li>Para cada componente principal, PCA encuentra un vector unitario centrado en el orígen apuntando en dirección del componente principal correspondiente.</li> <li>Recordemos que el vector unitario representa la dirección (positivo/negativo), y en nuestro caso es la dirección del PC en la que la varianza es la máxima. Los coeficientes (valores absoluto cada variable original son los "eigenvalues" que nos dicen que tanto contribuyen al cálculo del vector o PC.</li> <li># Con .components obtenemos los coeficientes o pesos de cada PC posComponents_df = pd.DataPrame(pcs.components_transpose(), columns = pcsSummary_df.columns, index = ndf.columns)</li> <li># Imprimimos los primeros 7 PC</li> </ul>
	<ul> <li>Examine la magnitud y la dirección de los coeficientes de las variables originales.</li> <li>Nota: Cuanto mayor sea el valor absoluto del coeficiente, más importante será la variable correspondiente en el cálculo del componente.</li> <li>Es importante mencionar que cada componente principal representa un vector con la combinación de las contribuciones de varianza de cada variable original.</li> <li>El PC1 captura la mayor cantidad de varianza (como se observa en la scree plot), sucesivamente el PC2 captura la mayor varianza posible bajo la condición de que este es ortogonal al es decir que no esta correlacionado con PC1. Los demás PCs siguen este mismo esquema.</li> <li>Para cada componente principal, PCA encuentra un vector unitario centrado en el orígen apuntando en dirección del componente principal correspondiente.</li> <li>Recordemos que el vector unitario representa la dirección (positivo/negativo), y en nuestro caso es la dirección del PC en la que la varianza es la máxima. Los coeficientes (valores absoluto cada variable original son los "eigenvalues" que nos dicen que tanto contribuyen al cálculo del vector o PC.</li> <li># Con .components obtenemos los coeficientes o pesos de cada PC possComponents_de possummary_df.columns, index = ndf.columns)</li> </ul>
	Examine la magnitud y la ciracción de los coeficientes de las variables originales.      Nota: Cuanto mayor see el valor absoluto del coeficiente, más importante será la variable correspondiente en el cálculo del componente.  Es importante menor y cada componente principal varianza (como se observa en la scree plot), sucesivamente el PC2 captura la mayor varianza posible bajo la condición de que este es ortogonal al es decir que no esta corribucioneo de conficientes (valores absoluto del componente principal). PCA encuentra un vector unitario respectado en el origen apuntando en dirección del componento principal correspondiente.  Recordemos que el vector unitario representa la dirección (positivo/hegativo), y en nuestro caso es la dirección del PC en la que la varianza e la máxima. Los coeficientes (valores absoluto del varianza e la máxima del varianza e la máxima la
	Examine la magnitud y la dirección de los coeficientes de las variables originales.  • Note: Cuamo moyor sea el valor absoluto del coeficiente, más importante será la variable correspondiente en el cálculo del componente.  Es importante mencionar: que cada componente principal representa un vector con la combinación de las contribuciones de varianta de cada variable original.  • El PCI captura la mayor cantidad de varianta (como se osserva en la sorse plotí, sucestivamente el PC2 captura la mayor varianza posible bajo la condición de que este es originales.  • Para cada componente principal. PCA encuentra un vector unitario centrado en el origina apuntando en clirección del componente principal correspondiente.  • Para cada componente principal. PCA encuentra un vector unitario centrado en el origina apuntando en clirección del componente principal correspondiente.  • Para cada componente principal. PCA encuentra un vector unitario centrado en el origina apuntando en clirección del componente principal correspondiente.  • Para cada componente principal correspondiente.  • Para cada componente principal correspondientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima. Los coeficientes (va orea absoluto cada varianza el la másima cada varianza el
	Examine la migritur y la ciracción de los conficientes de las variables originales.  Nota: Cuanto mayor ses el valor absoluto de conficiente, más importante será la variable correspondiente en el cáculo del componente.  Els importante membrane membrane que cartidad de incentionar que cartidad de incentionar que cartidad de interior como esta correspondiente en el cáculo del componente principal de cartidad en el valor de las contribuciones de varianza de cada variable original.  El PCC aquella a moyor cartidad de inventace Como as destructes de mayor cartidad de las contribuciones de varianza de cada variable original.  Para cada componente principal proprieta de membrane en la vez polity, acestivamente al PC2 captura la mayor varianza possible bipo a concisión de que ene es orrogenal al cada de componente principal
	Pose minor is magnified y a concessor of loc conficients de lac variables originals.      Pose des Cuarto magnified y a concessor de lacular de contractor principal representative de la variable componente.  Be importante mesclonar que cada componente principal representa un vector con la combinación de las contribuciones de varianza de cada variable original.  Be iPCC captura las responsables por cantidad de variables componente principal contractor de la cova piezo, successor de la cova piezo de componente principal de cartidad de la contribucione de varianza de cada variable original.  Pera cada componente principal por contractor de variables principal de cartidad de la contractor de la cova piezo de componente principal contractor de variables original.  Pera de cada componente principal por contractor de variables principal de cartidad de la contractor de la cova piezo de la cova piezo de variable principal contractor de variables de variables principal de cartidad de la contractor de la cova piezo de la cova piezo de variable principal contractor de variable principal contractor de variable de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de cartidad de la cova piezo de variable principal de variable principal de cartidad de la cova piezo de variable principal de variable p
	** Promised to Registed y to dissociate controllectors de las distance originales.  ** Biolita Courte progres de la cel de aboutée de controllectors de las distance de la controllectors de la controllector de la control
	Position in a manufaction of discounts of the position of
	Parameter   Para
	Process   Proc
	Process of the process of the control of the c
	The continue con
	Proposed a propos
	Part
	Continue
	Continue
	Control   Cont
	Company   Comp
	Part