

Actividad Semanal -- 4

Ciencia y analítica de Datos

Alumno: Alejandro Jesús Vázquez Navarro

Matrícula: A01793146

Parte 1

Ejercicio guiado

```
In [19]: 1 # Importo las librerías que necesitare para la extracción, transformación y limpieza de datos (si aplica)
          2
          3 import pandas as pd
          4 import numpy as np
          5
          6 # Agrego este valor para no truncar los resultados:
          7 pd.set_option('display.max_colwidth', None)
          8
```

```
In [20]: 1 # Realizo la adquisición de datos:
2 inPath = 'https://raw.githubusercontent.com/PosgradoMNA/Actividades_Aprendizaje-/main/default%20of%20cred
3 dfPCA = pd.read_csv(inPath, index_col = 0)
4
5 #Verifico que estén correctamente bajados
6
7 dfPCA.head()
```

Out[20]:

	X1	X2	X3	X4	X5	X6	X7	X8	X9	X10	...	X15	X16	X17	X18	X19	X20	X21	X22	X23
ID																				
1	20000	2.0	2.0	1.0	24.0	2.0	2.0	-1.0	-1.0	-2.0	...	0.0	0.0	0.0	0.0	689.0	0.0	0.0	0.0	0.0
2	120000	2.0	2.0	2.0	26.0	-1.0	2.0	0.0	0.0	0.0	...	3272.0	3455.0	3261.0	0.0	1000.0	1000.0	1000.0	0.0	2000.0
3	90000	2.0	2.0	2.0	34.0	0.0	0.0	0.0	0.0	0.0	...	14331.0	14948.0	15549.0	1518.0	1500.0	1000.0	1000.0	1000.0	5000.0
4	50000	2.0	2.0	1.0	37.0	0.0	0.0	0.0	0.0	0.0	...	28314.0	28959.0	29547.0	2000.0	2019.0	1200.0	1100.0	1069.0	1000.0
5	50000	1.0	2.0	1.0	57.0	-1.0	0.0	-1.0	0.0	0.0	...	20940.0	19146.0	19131.0	2000.0	36681.0	10000.0	9000.0	689.0	679.0

5 rows × 24 columns

In [21]:

```
1  # Realizo la asignación de cabeceras:
2
3  # Asignación de cabeceras de acuerdo a la información proporcionada:
4
5  #X1: Amount of the given credit (NT dollar): it includes both the individual consumer credit and his/her
6  #X2: Gender (1 = male; 2 = female).
7  #X3: Education (1 = graduate school; 2 = university; 3 = high school; 4 = others).
8  #X4: Marital status (1 = married; 2 = single; 3 = others).
9  #X5: Age (year).
10 #X6 - X11: History of past payment. We tracked the past monthly payment records (from April to September,
11 #X12-X17: Amount of bill statement (NT dollar). X12 = amount of bill statement in September, 2005; X13 =
12 #X18-X23: Amount of previous payment (NT dollar). X18 = amount paid in September, 2005; X19 = amount paid
13
14 # Creo un diccionario con los nombres de las columnas
15 col_dictionary = {'X1': 'given_credit',
16                  'X2': 'gender',
17                  'X3': 'education',
18                  'X4': 'marital_status',
19                  'X5': 'age',
20                  'X6': 'sep05_repayment',
21                  'X7': 'aug05_repayment',
22                  'X8': 'jul05_repayment',
23                  'X9': 'jun05_repayment',
24                  'X10': 'may05_repayment',
25                  'X11': 'abr05_repayment',
26                  'X12': 'sep05_bill_statement',
27                  'X13': 'aug05_bill_statement',
28                  'X14': 'jul05_bill_statement',
29                  'X15': 'jun05_bill_statement',
30                  'X16': 'may05_bill_statement',
31                  'X17': 'abr05_bill_statement',
32                  'X18': 'sep05_previous_payment',
33                  'X19': 'aug05_previous_payment',
34                  'X20': 'jul05_previous_payment',
35                  'X21': 'jun05_previous_payment',
36                  'X22': 'may05_previous_payment',
37                  'X23': 'abr05_previous_payment'
38                  }
39
40 # Renombro las columnas para tener un mejor entendimiento del dataset
41 dfPCA.rename(columns = col_dictionary , inplace=True)
```

```

42
43 dfPCA.head()
44 #df.columns
45

```

Out[21]:

	given_credit	gender	education	marital_status	age	sep05_repayment	aug05_repayment	jul05_repayment	jun05_repayment	may05_repayment
ID										
1	20000	2.0	2.0	1.0	24.0	2.0	2.0	-1.0	-1.0	-1.0
2	120000	2.0	2.0	2.0	26.0	-1.0	2.0	0.0	0.0	0.0
3	90000	2.0	2.0	2.0	34.0	0.0	0.0	0.0	0.0	0.0
4	50000	2.0	2.0	1.0	37.0	0.0	0.0	0.0	0.0	0.0
5	50000	1.0	2.0	1.0	57.0	-1.0	0.0	-1.0	-1.0	0.0

5 rows × 24 columns

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.

In [22]:

```

1 # Para este paso utilizaré matplotlib, numpy y PCA de sklearn decomposition
2
3 import matplotlib as mpl
4 import matplotlib.pyplot as plt
5 import numpy as np
6 from sklearn.decomposition import PCA
7
8 # Tengo 23 variables (incluyendo una dependiente Y que significa si el cliente cayó en default o no)

```

```
In [23]: 1 # Antes de comenzar debo eliminar las variables y para ello me basaré en el ejercicio anterior:
          2 # Esto es solo una sencilla validación:
          3
          4 dfPCA.isna().any()
          5
          6 # Tenemos nulos en casi todas las columnas excepto given_credit
          7
```

```
Out[23]: given_credit      False
          gender           True
          education        True
          marital_status    True
          age              True
          sep05_repayment   True
          aug05_repayment   True
          jul05_repayment   True
          jun05_repayment   True
          may05_repayment   True
          abr05_repayment   True
          sep05_bill_statement True
          aug05_bill_statement True
          jul05_bill_statement True
          jun05_bill_statement True
          may05_bill_statement True
          abr05_bill_statement True
          sep05_previous_payment True
          aug05_previous_payment True
          jul05_previous_payment True
          jun05_previous_payment True
          mayu05_previous_payment True
          abr05_previous_payment True
          Y                True
          dtype: bool
```

```
In [24]: 1 # Tengo variables categóricas como gender, education y marital_status, podría pasarles un one hot encodin
          2 # pero no es el objetivo de este ejercicio, así que las eliminaré del dataframe
          3
          4 dfOriginal = dfPCA
          5 dfPCA = dfPCA.loc[:, ~dfPCA.columns.isin(['gender', 'education', 'marital_status', 'Y'])]
          6
```

```
In [25]: 1 dfPCA.dropna(inplace=True)
2 print(f'Longitud del dataframe {len(dfPCA)}')
3
4 # El dataframe tiene 30k registros, se perdieron solamente 2.
5 # No haré una análisis de missing data porque esto fue parte de la tarea anterior
6
7
```

Longitud del dataframe 29958

C:\Users\aleja\AppData\Local\Temp\ipykernel_30632\2450569021.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy (https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)
dfPCA.dropna(inplace=True)

```
In [26]: 1 dfPCA.head()
```

Out[26]:

	given_credit	age	sep05_repayment	aug05_repayment	jul05_repayment	jun05_repayment	may05_repayment	abr05_repayment	sep
ID									
1	20000	24.0	2.0	2.0	-1.0	-1.0	-2.0	-2.0	
2	120000	26.0	-1.0	2.0	0.0	0.0	0.0	2.0	
3	90000	34.0	0.0	0.0	0.0	0.0	0.0	0.0	
4	50000	37.0	0.0	0.0	0.0	0.0	0.0	0.0	
5	50000	57.0	-1.0	0.0	-1.0	0.0	0.0	0.0	

In [27]:

```

1  # Comencemos con un tratamiento desde cero considerando todas las variables:
2  pcs = PCA(n_components=20)
3
4  pcs.fit(dfPCA)
5
6
7  pcsSummary = pd.DataFrame({'Standard deviation': np.sqrt(pcs.explained_variance_),
8                             'Proportion of variance': pcs.explained_variance_ratio_,
9                             'Cumulative proportion': np.cumsum(pcs.explained_variance_ratio_)}
10                             #column
11                             )
12  pcsSummary.round(4)

```

Out[27]:

	Standard deviation	Proportion of variance	Cumulative proportion
0	166585.0753	0.6106	0.6106
1	115824.7017	0.2952	0.9057
2	37256.8166	0.0305	0.9363
3	27747.1030	0.0169	0.9532
4	20698.6823	0.0094	0.9626
5	20278.0072	0.0090	0.9717
6	18520.7019	0.0075	0.9792
7	17040.2350	0.0064	0.9856
8	16294.1101	0.0058	0.9915
9	11615.2249	0.0030	0.9944
10	10413.1020	0.0024	0.9968
11	8786.7540	0.0017	0.9985
12	8197.6489	0.0015	1.0000
13	9.1140	0.0000	1.0000
14	2.0623	0.0000	1.0000
15	0.9386	0.0000	1.0000
16	0.7321	0.0000	1.0000

	Standard deviation	Proportion of variance	Cumulative proportion
17	0.5963	0.0000	1.0000
18	0.5120	0.0000	1.0000

In [28]:

```

1 # Hagamos una transposición para mejor lectura
2
3 pcsSummary = pcsSummary.transpose()
4 pcsSummary.round(2)

```

Out[28]:

	0	1	2	3	4	5	6	7	8	9	10	11	
Standard deviation	166585.08	115824.70	37256.82	27747.10	20698.68	20278.01	18520.70	17040.23	16294.11	11615.22	10413.1	8786.75	819
Proportion of variance	0.61	0.30	0.03	0.02	0.01	0.01	0.01	0.01	0.01	0.00	0.0	0.00	
Cumulative proportion	0.61	0.91	0.94	0.95	0.96	0.97	0.98	0.99	0.99	0.99	1.0	1.00	

```

1 Con esta tabla podemos calcular fácilmente que:
2
3 * `z1` accounts for 61.01% del total de varianza y
4 * `z2` aporta un total de 29.51%
5
6 Esto suma un 90.57% del total de varianza.
7
8 Significa pues que si removemos el PC2 tendríamos aún 61.01% de varianza. Este componente corresponde al
  crédito otorgado
9

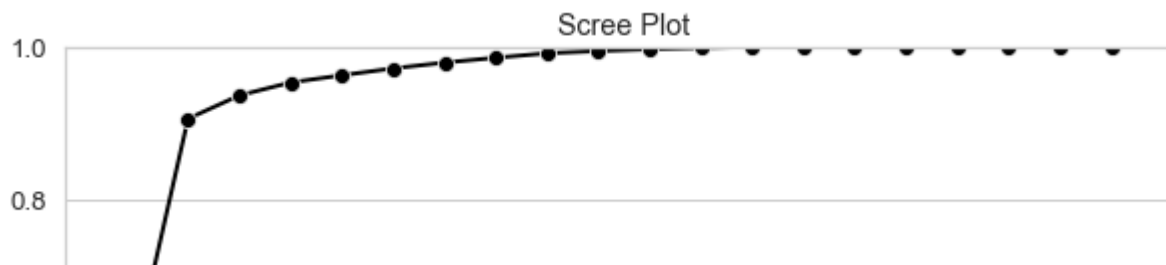
```

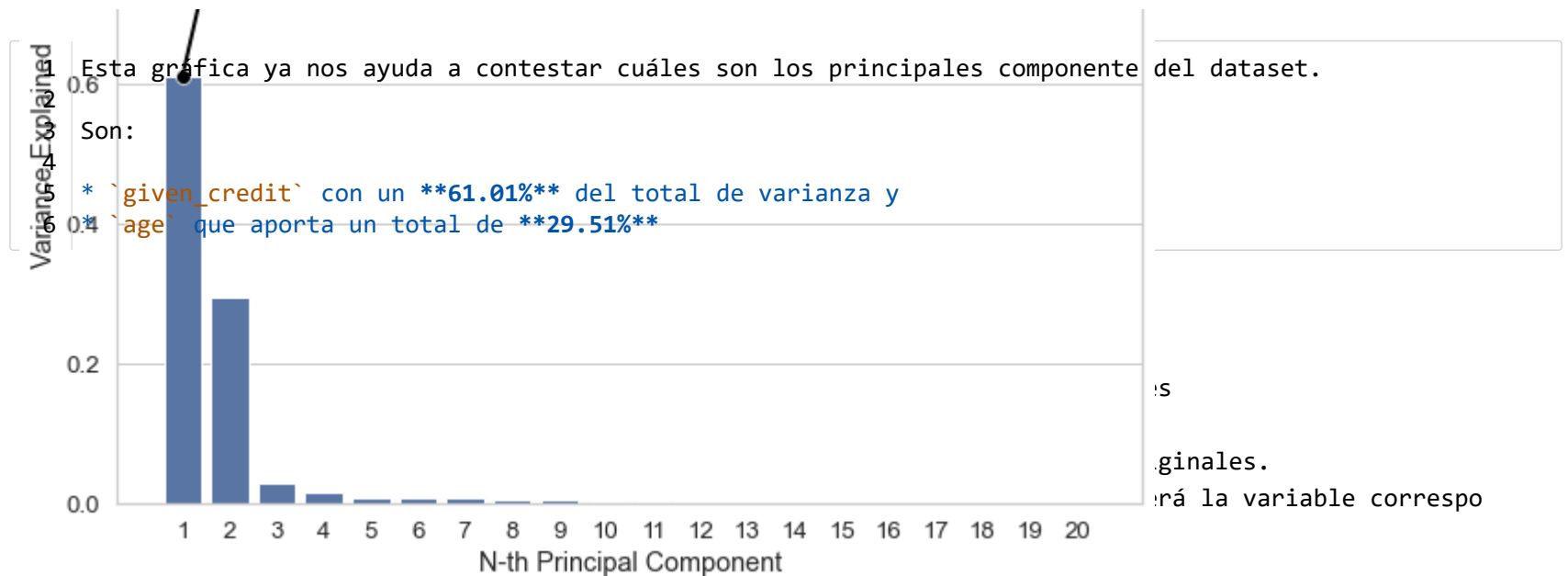

In [29]:

```

1  # Hagamos una gráfica para representarlo y determinar el número mínimo de componentes:
2
3  import seaborn as sns
4
5  PC_components = np.arange(pcs.n_components_) + 1
6  #PC_components
7
8  _ = sns.set(style = 'whitegrid',
9              font_scale = 1.2
10             )
11
12  fig, ax = plt.subplots(figsize=(10, 7))
13
14  _ = sns.barplot(x = PC_components,
15                 y = pcs.explained_variance_ratio_,
16                 color = 'b'
17                 )
18
19  _ = sns.lineplot(x = PC_components-1,
20                  y = np.cumsum(pcs.explained_variance_ratio_),
21                  color = 'black',
22                  linestyle = '-',
23                  linewidth = 2,
24                  marker = 'o',
25                  markersize = 8
26                  )
27
28  plt.title('Scree Plot')
29  plt.xlabel('N-th Principal Component')
30  plt.ylabel('Variance Explained')
31  plt.ylim(0, 1)
32  plt.show()
33

```





```
In [30]: 1 # Examinemos los componentes:
          2
          3
          4 pcsSummary.columns = ['PC{}'.format(i) for i in range(1, len(pcsSummary.columns) + 1)]
          5 pcsSummary.round(4)
```

Out[30]:

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10
Standard deviation	166585.0753	115824.7017	37256.8166	27747.1030	20698.6823	20278.0072	18520.7019	17040.2350	16294.1101	11615.2249
Proportion of variance	0.6106	0.2952	0.0305	0.0169	0.0094	0.0090	0.0075	0.0064	0.0058	0.0030
Cumulative proportion	0.6106	0.9057	0.9363	0.9532	0.9626	0.9717	0.9792	0.9856	0.9915	0.9944

```

In [31]: 1 # Es necesario identificar magnitudes y direcciones de cada componente para determinar importancia de cad
          2 # en cada componente:
          3
          4 pcsComponents_dfX = pd.DataFrame(pcs.components_.transpose(),
          5                                   columns = pcsSummary.columns,
          6                                   index = dfPCA.iloc[:, 0:].columns
          7                                   )
          8
          9 # Redondeo a 4 cifras para una mejor lectura
         10 pcsComponents_dfX.round(4)

```

Out[31]:

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10	PC11	PC12	PC13
given_credit	0.4911	0.8693	-0.0206	-0.0184	-0.0441	0.0169	-0.0011	0.0010	0.0095	0.0038	0.0058	0.0000	0.0019
age	0.0000	0.0000	-0.0000	-0.0000	-0.0000	0.0000	-0.0000	-0.0000	0.0000	-0.0000	0.0000	-0.0000	-0.0000
sep05_repayment	0.0000	-0.0000	-0.0000	-0.0000	-0.0000	0.0000	-0.0000	0.0000	-0.0000	0.0000	0.0000	0.0000	-0.0000
aug05_repayment	0.0000	-0.0000	-0.0000	-0.0000	-0.0000	0.0000	-0.0000	0.0000	-0.0000	0.0000	0.0000	0.0000	0.0000
jul05_repayment	0.0000	-0.0000	0.0000	-0.0000	-0.0000	0.0000	0.0000	0.0000	-0.0000	0.0000	0.0000	0.0000	0.0000
jun05_repayment	0.0000	-0.0000	0.0000	0.0000	-0.0000	0.0000	0.0000	0.0000	-0.0000	0.0000	-0.0000	0.0000	0.0000
may05_repayment	0.0000	-0.0000	0.0000	-0.0000	-0.0000	-0.0000	-0.0000	-0.0000	-0.0000	0.0000	0.0000	0.0000	-0.0000
abr05_repayment	0.0000	-0.0000	0.0000	-0.0000	-0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	-0.0000	0.0000
sep05_bill_statement	0.3886	-0.2211	-0.5489	-0.2049	0.3169	0.0026	-0.4563	0.0073	0.0703	0.1727	-0.2674	-0.2046	-0.0059
aug05_bill_statement	0.3815	-0.2262	-0.3910	-0.0378	0.0115	-0.0805	0.5269	-0.0660	0.0993	0.0079	0.3511	0.4765	0.0160
jul05_bill_statement	0.3723	-0.2163	-0.0481	0.5534	-0.2476	0.2958	0.0287	0.0813	-0.1080	-0.3365	0.1515	-0.4416	-0.0825
jun05_bill_statement	0.3465	-0.1938	0.2610	0.0791	-0.3679	-0.4732	-0.1782	-0.2277	-0.0974	-0.1369	-0.3749	0.2335	0.3189
may05_bill_statement	0.3230	-0.1766	0.4281	-0.2001	-0.0158	-0.0783	-0.0134	0.4801	0.1265	0.1010	-0.0525	0.0924	-0.6059
abr05_bill_statement	0.3087	-0.1672	0.4837	-0.3084	0.2269	0.3918	0.0768	-0.2629	-0.1055	0.2368	0.2031	-0.1585	0.3662
sep05_previous_payment	0.0266	0.0057	0.0375	0.1935	0.2025	-0.2280	0.6117	-0.0826	0.1709	0.2792	-0.4709	-0.3965	-0.0172
aug05_previous_payment	0.0313	0.0108	0.1601	0.6733	0.3373	0.1076	-0.2548	-0.0047	0.2264	0.3581	0.0235	0.3840	0.0470
jul05_previous_payment	0.0268	0.0110	0.1351	0.0283	0.3106	-0.5666	-0.1526	-0.3082	0.2279	-0.2016	0.4972	-0.2858	-0.1512
jun05_previous_payment	0.0222	0.0104	0.0944	-0.0403	0.4364	0.0203	0.0854	0.4481	0.2273	-0.5907	-0.1465	0.0790	0.4020

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9	PC10	PC11	PC12	PC13
may05_previous_payment	0.0222	0.0117	0.0550	0.0137	0.2551	0.2087	0.0602	0.4816	0.2075	0.2028	0.2000	0.2246	0.4446

```

1 Con esta tabla podemos interpretar fácilmente lo siguiente:
2
3 - El primer componente principal (PC1) está dominado por dos variables `given_credit` con una magnitud
   positiva de 0.4911 y por `sep05_bill_statement` con un peso de 0.3886
4
5 - El segundo componente está completamente dominado por la variable `given_credit` con un peso de
   0.8693
6
7 - No estoy agregando el PC3 debido a su poca aportación al total de varianza acumulada: 3.05

```

Paso 3:

- Identifique valores atípicos
- Realice alguna gráfica de valores atípicos o boxplot para identificar los valores atípicos.
- Cualquier punto que esté más alejado de la línea de referencia es un valor atípico.

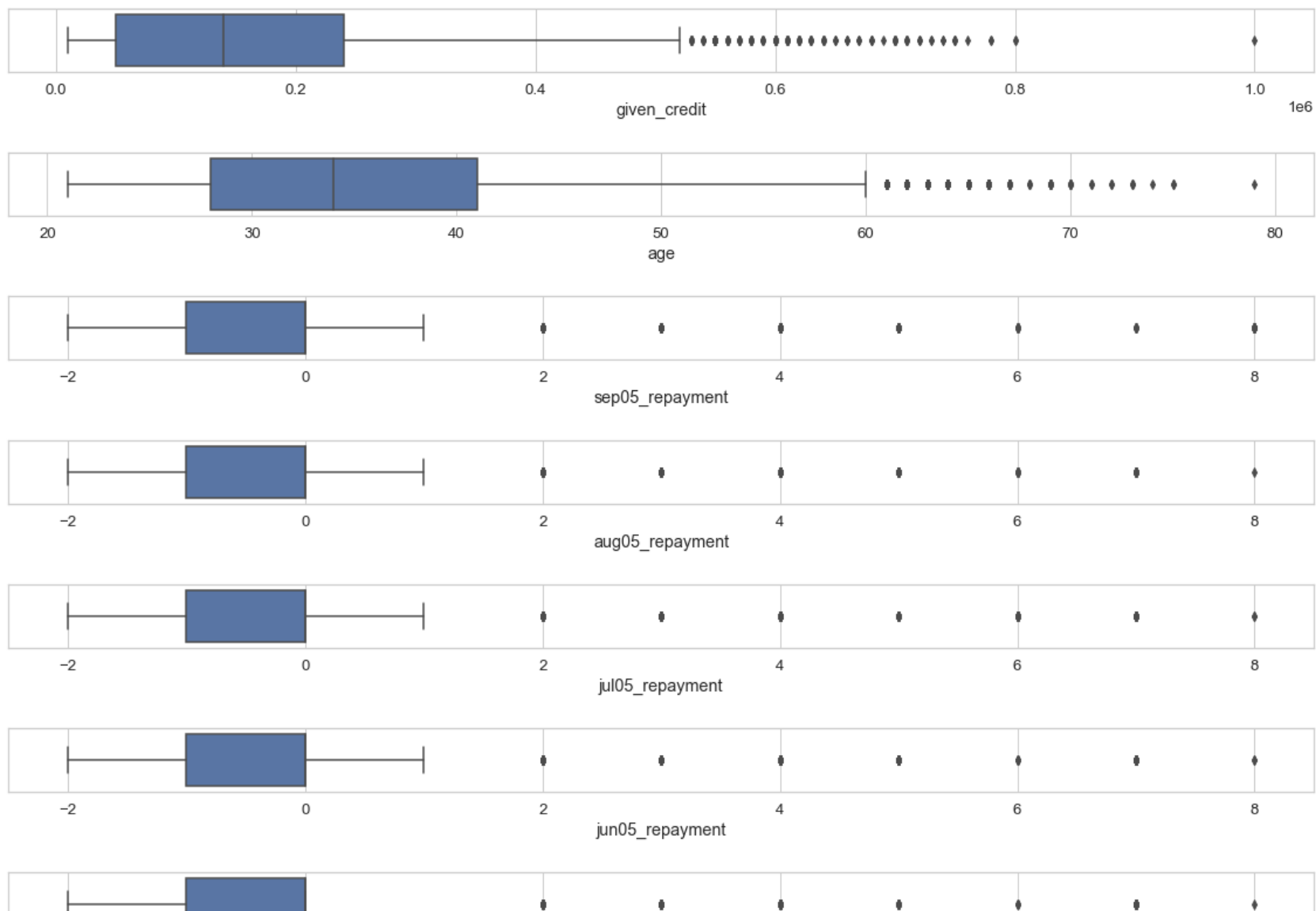
```
In [32]: 1 # Para este cometido utilizaré el dataframe generado : pcsComponents_dfX
2
3 pcsComponents_dfX['feature'] = pcsComponents_dfX.index
4
5 dfOutLiers=pcsComponents_dfX.reset_index()
6 dfOutLiers.reset_index(inplace=True)
7
8 # Elimino una variable que surgió
9 dfOutLiers.pop("level_0")
10 dfOutLiers.pop("index")
11
12
13 # Verifico mi dataframe
14 dfOutLiers.head()
15 dfOutLiers.columns
16
```

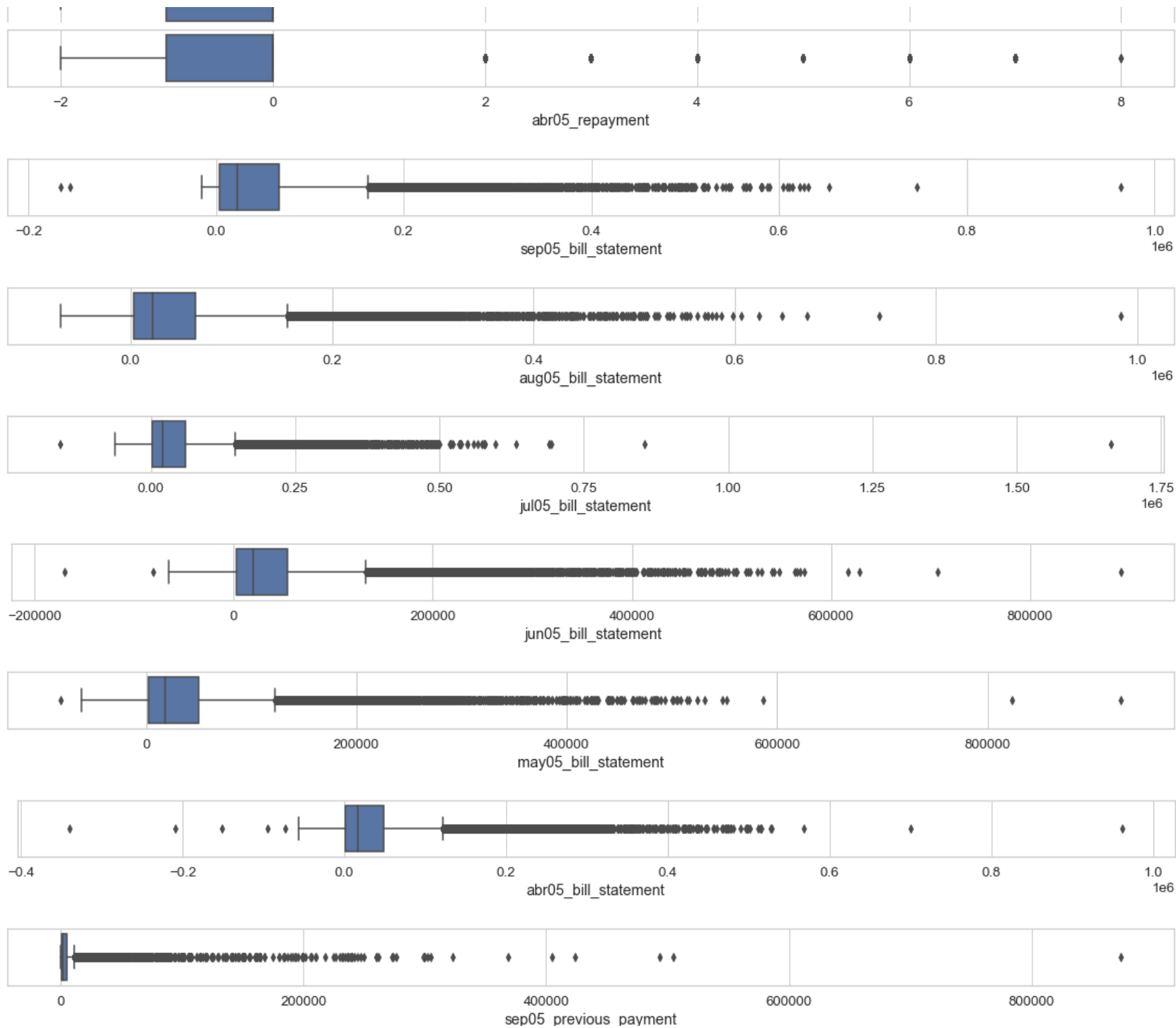
```
Out[32]: Index(['PC1', 'PC2', 'PC3', 'PC4', 'PC5', 'PC6', 'PC7', 'PC8', 'PC9', 'PC10',
               'PC11', 'PC12', 'PC13', 'PC14', 'PC15', 'PC16', 'PC17', 'PC18', 'PC19',
               'PC20', 'feature'],
              dtype='object')
```

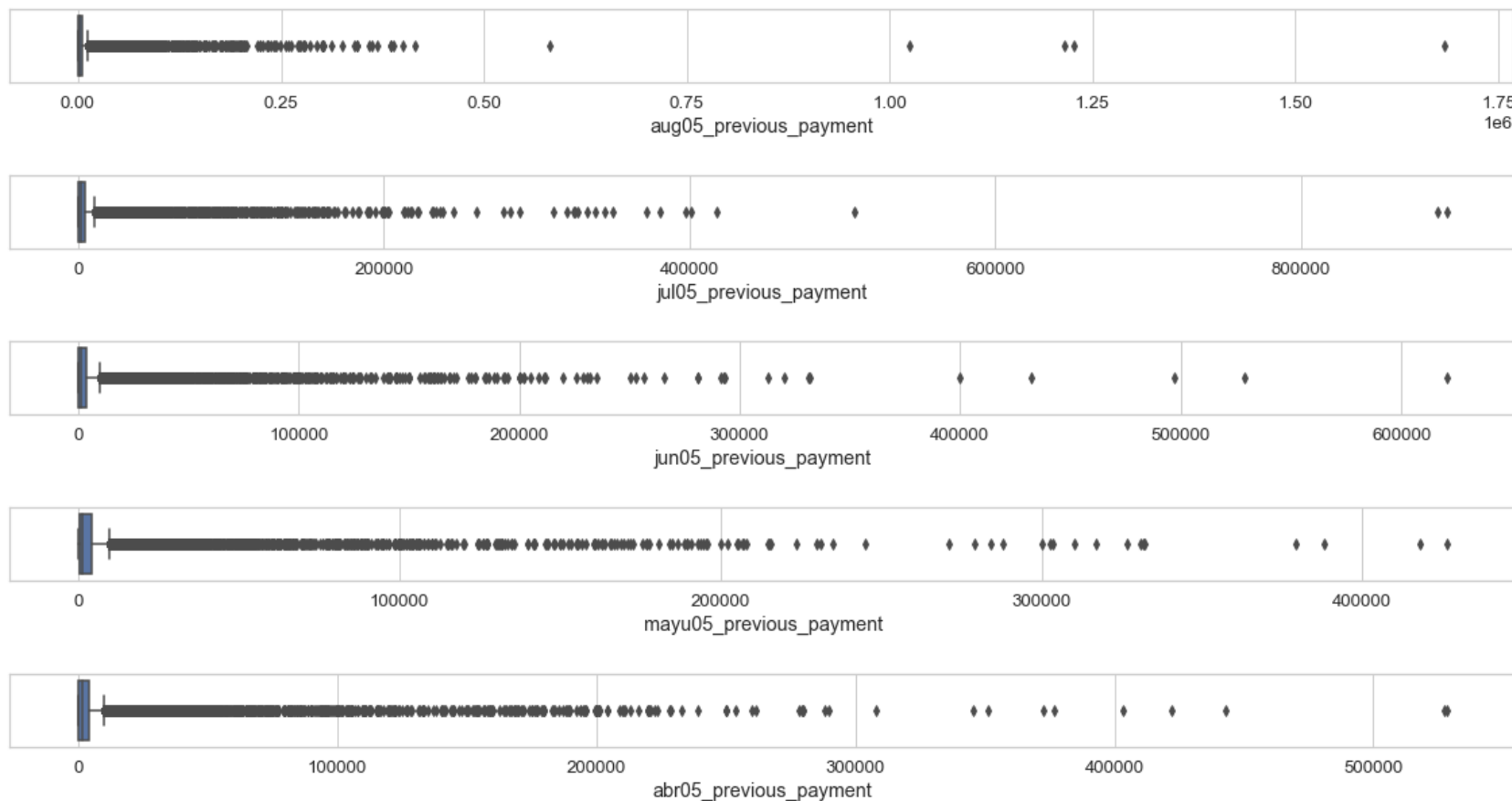
```
In [44]: 1 # Esta línea es solo cosmética para no hacer scroll down
2
3 %%javascript
4 IPython.OutputArea.prototype._should_scroll = function(lines) {
5     return false;
6 }
7
```

In [45]:

```
1 plt.rcParams["figure.figsize"] = (12,10)
2
3
4
5 for column in dfPCA:
6     plt.figure(figsize=(20,1))
7     sns.boxplot(data=dfOriginal, x=column)
```







Interesante porque podemos encontrar outliers en `given credit`, `age`, `sep05_bill_statement`, `aug05_bill_statement`, `jun05_bill_statement`, `may05_bill_statement`, y en general en todas las columnas de pago previo

Parte 2

Responde las siguientes preguntas en una celda de texto en Jupyter Notebook

- ¿Cuál es el número de componentes mínimo y por qué?

Como el análisis de varianza fue posible encontrar que:

- z1 aporta el 61.01% del total de varianza y
- z2 aporta un total de 29.51%

Esto suma un 90.57% del total de varianza. No es necesario incluir los demás componentes pues tienen una aportación muy baja.

Significa pues que si removemos el PC2 tendríamos aún 61.01% de varianza. Este componente corresponde al crédito otorgado.

Estos componentes son:

- given_credit
- age

- **¿Cuál es la variación de los datos que representan esos componentes?**

Con el cálculo de varianza se pudo calcular que:

- z1 aporta el 61.01% del total de varianza y
- z2 aporta un total de 29.51%

- **¿Cuál es la pérdida de información después de realizar PCA?**

Podemos quedarnos con el 90.57% de la información y perder 9.43%

- **De las variables originales, ¿Cuál tiene mayor y cuál tiene menor importancia en los componentes principales?**

Después de calcular el peso de cada variable para cada componente principal pude calcular que:

- El primer componente principal (PC1) está dominado por dos variables given_credit con una magnitud positiva de 0.4911 y por sep05_bill_statement con un peso de **0.3886**
- El segundo componente está completamente dominado por la variable given_credit con un peso de **0.8693**
- No estoy agregando el PC3 debido a su poca aportación al total de varianza acumulada: **3.05**

- **¿Cuándo se recomienda realizar un PCA y qué beneficios ofrece para Machine Learning?**

La reducción de dimensionalidad es una técnica necesaria cuando tenemos un número alto de características en nuestro

dataset. Además de tener un nivel alto de correlación entre diferentes variables, eso lo podemos conocer como multicolinealidad. PCA aplica un principio de parsimonia donde obtenemos el máximo de información con el mínimo de características.

Como beneficios principales, nos ofrece la reduccion de la matriz del dataset reteniendo la máxima información posible, esto deriva en lo siguiente:

- Entendimiento más sencillo de las variables que aportan información al modelo.
- Reducción de procesamiento pues tendremos matrices más pequeñas.
- Naturalmente ayuda a reducir el overfitting pues nuestro modelo tomará solamente las variables necesarias.