

Analisis de Clientes Banco ACME S.A

Customer Behavior de Tarjetas de Credito

Objetivos del analisis

Este Notebook tiene los siguientes objetivos:

- Explorar de base de datos de clientes, utilizando varios tipos de Visualización.
- Realizar el preprocesado de datos, preparandolos para el analisis y modelado.
- Agrupar clientes utilizando modelos de clusterización.
- Interpretar analisis sobre los grupos de clientes, perfilado de los grupos.
- Elevar propuesta de marketing en base a los perfiles y análisis realizados.

Librerias Utilizadas

In [269]:

```
1 !conda activate credit_cards
```

In [270]:

```
1 # --- Importing Libraries ---
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import pandas_profiling
6 #from ydata_profiling import ProfileReport
7
8 import scipy.stats as stats
9
10 #import ydata_profiling
11
12 import seaborn as sns
13 import warnings
14 import os
15 import yellowbrick
16 import scipy.cluster.hierarchy as shc
17 import matplotlib.patches as patches
18
19 from faker import Faker
20 import random
21
22 from ipywidgets import interact, interactive, fixed, interact_manual
23 import ipywidgets as widgets
24 from IPython import display
25
26
27 from matplotlib.patches import Rectangle
28 from pandas_profiling import ProfileReport
29 from pywaffle import Waffle
30 from math import isnan
31 from random import sample
32 from numpy.random import uniform
33 from sklearn.neighbors import NearestNeighbors
34 from sklearn.impute import KNNImputer
35 from sklearn.preprocessing import StandardScaler
36 from sklearn.decomposition import PCA
37 from sklearn.cluster import KMeans, DBSCAN, AgglomerativeClustering
38 from sklearn.metrics import davies_bouldin_score, silhouette_score, calinski_harabasz
39 from yellowbrick.cluster import KElbowVisualizer, SilhouetteVisualizer
40 from yellowbrick.style import set_palette
41 from yellowbrick.contrib.wrapper import wrap
42
43 # --- Libraries Settings ---
44 warnings.filterwarnings('ignore')
45 sns.set_style('white')
46 plt.rcParams['figure.dpi'] = 600
47 #sns.set(rc = {'axes.facecolor': '#FBFBFB', 'figure.facecolor': '#FBFBFB'})
48 class clr:
49     start = '\033[93m'+'\033[1m'
50     color = '\033[93m'
51     end = '\033[0m'
```

Creación del Dataset a utilizar

In [271]:

```
1 df = pd.read_csv('Customer_Data.csv')
2
3 # --- Reading Train Dataset ---
4 print(clr.start+'.: Dataset Importado :.'+clr.end)
5 print(cclr.color+'*' * 23)
6 data.head().style.background_gradient(cmap='YlOrBr')
```

```
.: Dataset Importado :.
*****
```

Out[271]:

	cust_id	balance	balance_frequency	purchases	oneoff_purchases	installments_purc
0	C10001	40.900749	0.818182	95.400000	0.000000	95.4
1	C10002	3202.467416	0.909091	0.000000	0.000000	0.0
2	C10003	2495.148862	1.000000	773.170000	773.170000	0.0
3	C10004	1666.670542	0.636364	1499.000000	1499.000000	0.0
4	C10005	817.714335	1.000000	16.000000	16.000000	0.0

Descripcion de los datos

- Balance: Es el monto que el cliente posee disponible para seguir gastando
- Balance Frequency: Valor entre 0 y 1, determinar que tan frecuente es el movimiento de la tarjeta
- Purchases: Monto total invertido en compras
- oneoff_purchases: El monto total invertido en compras al contado
- Installments_purchases: Es el monto invertido en compras a cuotas
- Cash_advance: Monto que el clientes saca en cajeros / Adelantos de efectivo
- Purchases_frequency: Frecuencia de compras
- Oneoff_purchases_rfecency: frecuencia con la que el cliente realiza compras al contado
- Purchases_installments_frequency: frecuencia con la que el cliente realiza compras a cuotas
- Cash_advance_frequency: Frecuencia con la que el cliente realiza adelantos de efectivo
- Cash_advance_trx: Cantidad de transacciones de adelanto de efectivo
- purchases_trx: Cantidad total de transacciones
- credit_limit: Limite de credito del cliente
- payments: Montos ya pagados sobre la tarjeta
- minimum_payments:
- prc_full_payment:
- tenure:

Se adicionaran los siguientes datos:

- Antigüedad del cliente (En meses)
- Antigüedad del cliente con el productos (En Meses)
- Ciudad
- Edad
- Genero

- Rubro de mayor consumo (Categoria)

Seria importante agregar:

- Es digital? (Podria ser "Frecuencia de transacciones digitales")
- Es cliente final de "Pago de salario"
- Cantidad de productos del cliente
- Rubro de mayor consumo (Categoria)
- TMHI: Tiempo minimo hasta la inactiviad
- Promedio de ingresos del cliente

In [272]:

```
1 df.head()
```

Out[272]:

	cust_id	balance	balance_frequency	purchases	oneoff_purchases	installments_purcha
0	C10001	40.900749	0.818182	95.40	0.00	
1	C10002	3202.467416	0.909091	0.00	0.00	
2	C10003	2495.148862	1.000000	773.17	773.17	
3	C10004	1666.670542	0.636364	1499.00	1499.00	
4	C10005	817.714335	1.000000	16.00	16.00	

Revision Pre-exploracion del dataset

In [273]:

```
1 # --- Dataset Report ---
2 ProfileReport(df, title='Credit Card Dataset Report', minimal=True, progress_bar=False)
```

In [274]:

```
1 # eliminar las columnas 'cust_id', "ciudad", "genero", "Top_Rubro" por que no son numerici
2 df = df.drop(columns=['cust_id'])
```

In [275]:

```
1 df.head()
```

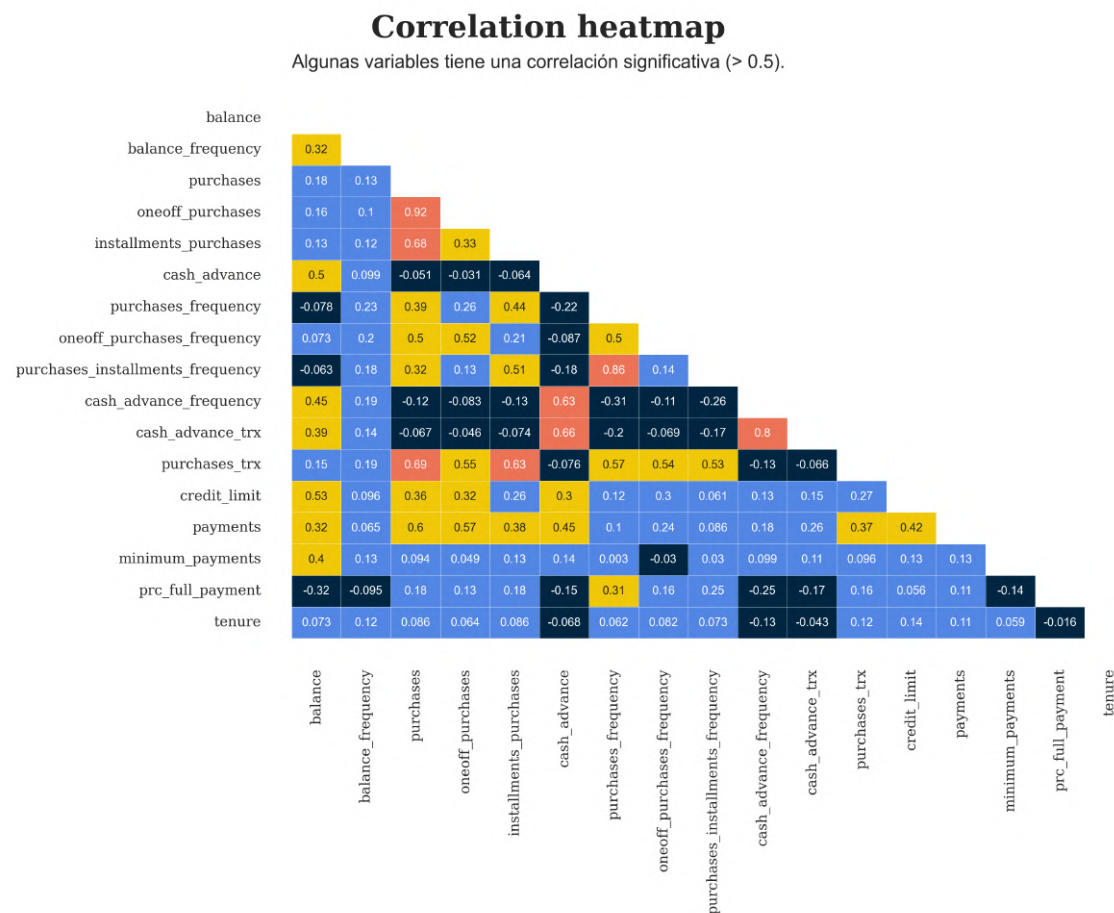
Out[275]:

	balance	balance_frequency	purchases	oneoff_purchases	installments_purchases	cas
0	40.900749	0.818182	95.40	0.00	95.4	
1	3202.467416	0.909091	0.00	0.00	0.0	6
2	2495.148862	1.000000	773.17	773.17	0.0	
3	1666.670542	0.636364	1499.00	1499.00	0.0	
4	817.714335	1.000000	16.00	16.00	0.0	

Correlación entre variables

In [276]:

```
1 %matplotlib inline
2 # --- Correlation Map (Heatmap) ---
3 mask = np.triu(np.ones_like(df.corr(), dtype=bool))
4 fig, ax = plt.subplots(figsize=(7, 6))
5 sns.heatmap(df.corr(), mask=mask, annot=True, cmap= ['#002642', '#5386E4', '#F0C808'],
6 yticks, ylabels = plt.yticks()
7 xticks, xlabels = plt.xticks()
8 ax.set_xticklabels(xlabels, size=6, fontfamily='serif')
9 ax.set_yticklabels(ylabels, size=6, fontfamily='serif')
10 plt.suptitle('Correlation heatmap', fontweight='heavy', x=0.327, y=0.96, ha='left', f
11 plt.title('Algunas variables tiene una correlación significativa (> 0.5).\n', fontsize
12 plt.tight_layout(rect=[0, 0.04, 1, 1.01])
13 plt.show(block=True)
```



EDA

In [277]:

```
1 df.columns = df.columns.str.upper()
```

EDA 1: Balance del cliente vs Limite de Credito

In [278]:

```

1 # --- EDA 1 Variables ---
2 scatter_style=dict(linewidth=0.65, edgecolor='#100C07', alpha=0.85)
3 sub_scatter_style_color=dict(s=5, alpha=0.65, linewidth=0.15, zorder=10, edgecolor='#
4 sub_scatter_style_grey=dict(s=5, alpha=0.3, linewidth=0.7, zorder=5, color='#CAC9CD')
5 grid_style=dict(alpha=0.3, color='#9B9A9C', linestyle='dotted', zorder=1)
6 xy_label=dict(fontweight='bold', fontsize=14, fontfamily='serif')
7 suptitle=dict(fontsize=22, fontweight='heavy', fontfamily='serif')
8 title=dict(fontsize=16, fontfamily='serif')
9 color_pallete=['#2D0F51', '#FF9A00', '#6600A5', '#FFD61E', '#722E9A', '#FFE863', '#A4
10 sub_axes=[None] * 7
11
12 # --- EDA 1 Data Frame ---
13 eda1 = df[['CREDIT_LIMIT', 'BALANCE', 'TENURE']]
14 eda1['TENURE'] = eda1['TENURE'].astype(str)
15 tenure = sorted(eda1['TENURE'].unique())
16
17 # --- EDA 1 Settings ---
18 fig = plt.figure(figsize=(22, 14))
19 gs = fig.add_gridspec(7, 7)
20 ax = fig.add_subplot(gs[:, :7])
21 ax.set_aspect(1)
22
23 # --- EDA 1: Main Scatter Plot ---
24 for x in range(len(tenure)):
25     eda1_x = eda1[eda1['TENURE']==tenure[x]]
26     ax.scatter(eda1_x['CREDIT_LIMIT'], eda1_x['BALANCE'], s=80, color=color_pallete[x]
27     ax.set_title('There are positive correlation between both variables. Most credit
28     ax.set_xlabel('\nCREDIT_LIMIT', **xy_label)
29     ax.set_ylabel('BALANCE\n', **xy_label)
30     ax.grid(axis='y', which='major', **grid_style)
31     ax.grid(axis='x', which='major', **grid_style)
32     for spine in ax.spines.values():
33         spine.set_color('None')
34     for spine in ['bottom', 'left']:
35         ax.spines[spine].set_visible(True)
36         ax.spines[spine].set_color('#CAC9CD')
37     plt.xticks(fontsize=12)
38     plt.yticks(fontsize=12)
39
40 # --- EDA 1: Sub Plots ---
41 for idx, tnr in enumerate(tenure):
42     sub_axes[idx] = fig.add_subplot(gs[idx, 6], aspect=1)
43
44     sub_axes[idx].scatter(eda1[eda1['TENURE']!=tnr]['CREDIT_LIMIT'], eda1[eda1['TENUR
45     sub_axes[idx].scatter(eda1[eda1['TENURE']==tnr]['CREDIT_LIMIT'], eda1[eda1['TENUR
46
47     cnt = (eda1['TENURE']==tnr).sum()
48     sub_axes[idx].set_title(f'Tenure {tnr} - ({cnt})', loc='left', fontsize=10, fontf
49     sub_axes[idx].set_xticks([])
50     sub_axes[idx].set_yticks([])
51     for spine in sub_axes[idx].spines.values():
52         spine.set_color('None')
53
54 # --- EDA 1 XY Limit ---
55 for axes in [ax] + sub_axes:
56     axes.set_xlim(-1000, 31000)
57     axes.set_ylim(-1000, 20000)
58
59 # --- EDA 1 Title ---

```

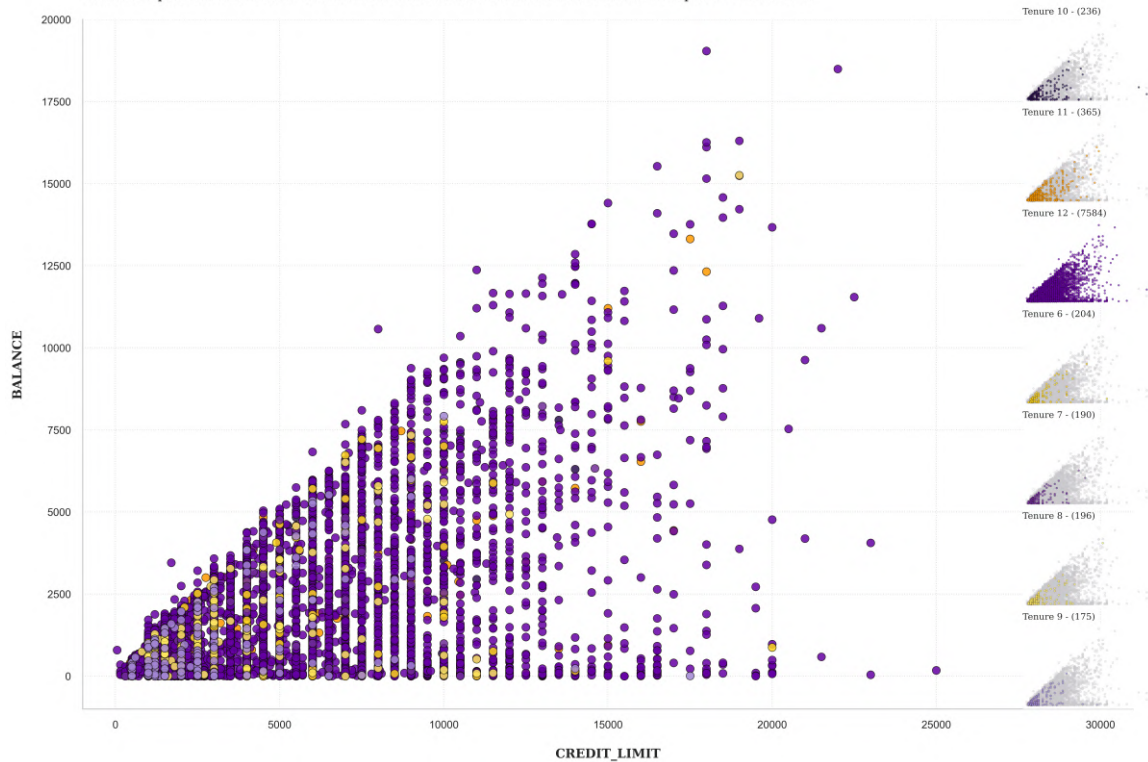
```

60 plt.suptitle('Scatter Plot Credit Limit vs. Balance based on Tenure', x=0.138, y=0.94
61
62 plt.show());

```

Scatter Plot Credit Limit vs. Balance based on Tenure

There are positive correlation between both variables. Most credit card customers prefer 12 months.



EDA 2: Promedio, Minimo y Maximo de "Monto de Compra" y "Cantidad de Transacciones"

Agrupado por Tenure del cliente

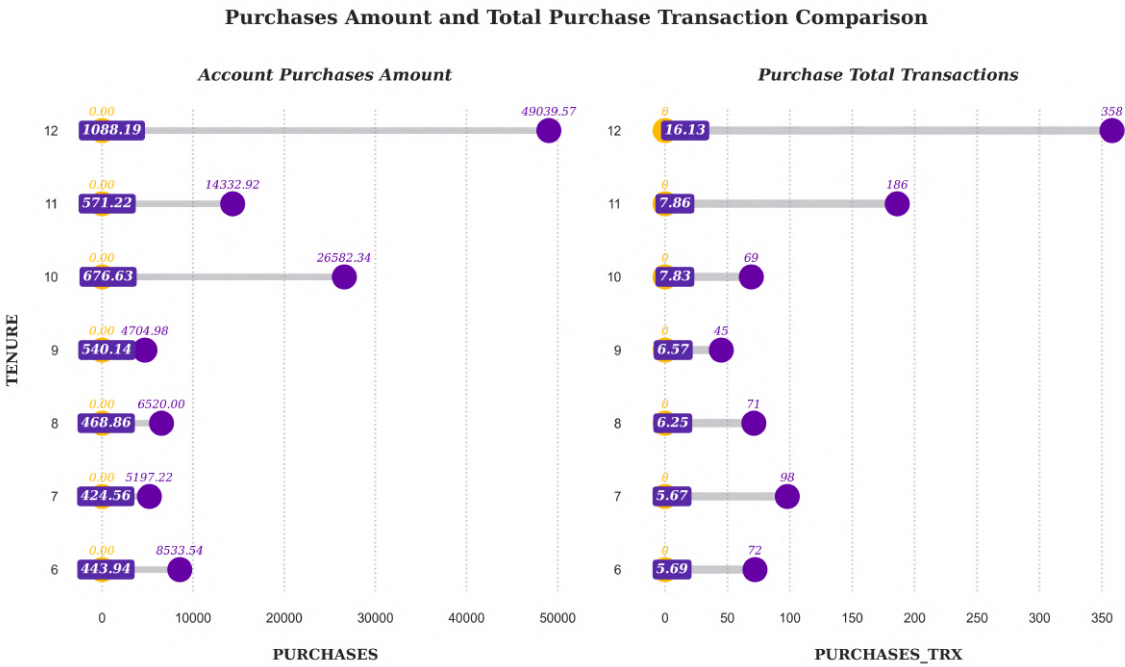
In [279]:

```

1 # --- EDA 2 Variables ---
2 title=dict(fontsize=10, fontfamily='serif', style='italic', weight='bold', ha='center')
3 grid_style = dict(alpha=0.6, color='#9B9A9C', linestyle='dotted', zorder=1)
4 sct_style = dict(s=175, linewidth=2)
5 xy_label = dict(fontweight='bold', fontsize=9, fontfamily='serif')
6 ann_style = dict(xytext=(0, 0), textcoords='offset points', va='center', ha='center',
7 tenure = sorted(df['TENURE'].unique())
8 color_pallete = ['#2D0F51', '#FF9A00', '#6600A5', '#FFD61E', '#722E9A', '#FFE863', '#
9
10 # --- EDA 2.1 Data Frame ---
11 eda2_1 = df[['PURCHASES', 'TENURE']]
12 eda2_1 = eda2_1.groupby('TENURE').agg(MIN=('PURCHASES', 'min'), AVG=('PURCHASES', 'me
13
14 # --- EDA 2.2 Data Frame ---
15 eda2_2 = df[['PURCHASES_TRX', 'TENURE']]
16 eda2_2 = eda2_2.groupby('TENURE').agg(MIN=('PURCHASES_TRX', 'min'), AVG=('PURCHASES_T
17
18 # --- EDA 2.1 & 2.2 Settings ---
19 fig = plt.figure(figsize=(10, 6))
20 plt.suptitle('\nPurchases Amount and Total Purchase Transaction Comparison', fontweig
21
22 # --- EDA 2.1 (Left Dumbbell) ---
23 plt.subplot(1, 2, 1)
24 plt.tight_layout(rect=[0, 0, 1, 1.01])
25 axs_left=plt.gca()
26 min_sct = plt.scatter(x=eda2_1['MIN'], y=eda2_1['TENURE'], c='#FFBB00', **sct_style)
27 max_sct = plt.scatter(x=eda2_1['MAX'], y=eda2_1['TENURE'], c='#6600A5', **sct_style)
28 for i in range(len(tenure)):
29     eda2_1_x = eda2_1[eda2_1['TENURE']==tenure[i]]
30     plt.hlines(y=eda2_1_x['TENURE'], xmin=eda2_1_x['MIN'], xmax=eda2_1_x['MAX'], line
31     plt.annotate('{0:.2f}'.format(eda2_1_x['MIN'].values[0]), xy=(eda2_1_x['MIN'].valu
32     plt.annotate('{0:.2f}'.format(eda2_1_x['AVG'].values[0]), xy=(eda2_1_x['AVG'].valu
33     plt.annotate('{0:.2f}'.format(eda2_1_x['MAX'].values[0]), xy=(eda2_1_x['MAX'].valu
34 for spine in axs_left.spines.values():
35     spine.set_color('None')
36 plt.xlabel('\nPURCHASES', **xy_label)
37 plt.ylabel('TENURE\n', **xy_label)
38 plt.xticks(fontsize=8)
39 plt.yticks(fontsize=8)
40 plt.grid(axis='y', alpha=0)
41 plt.grid(axis='x', which='major', **grid_style)
42 plt.title('\nAccount Purchases Amount\n', **title)
43
44 # --- EDA 2.2 (Right Dumbbell) ---
45 plt.subplot(1, 2, 2)
46 plt.tight_layout(rect=[0, 0, 1, 1.01])
47 axs_right=plt.gca()
48 min_sctt = plt.scatter(x=eda2_2['MIN'], y=eda2_2['TENURE'], c='#FFBB00', **sct_style)
49 max_sctt = plt.scatter(x=eda2_2['MAX'], y=eda2_2['TENURE'], c='#6600A5', **sct_style)
50 for i in range(len(tenure)):
51     eda2_2_x = eda2_2[eda2_2['TENURE']==tenure[i]]
52     plt.hlines(y=eda2_2_x['TENURE'], xmin=eda2_2_x['MIN'], xmax=eda2_2_x['MAX'], line
53     plt.annotate('{:.0f}'.format(eda2_2_x['MIN'].values[0]), xy=(eda2_2_x['MIN'].valu
54     plt.annotate('{0:.2f}'.format(eda2_2_x['AVG'].values[0]), xy=(eda2_2_x['AVG'].valu
55     plt.annotate('{:.0f}'.format(eda2_2_x['MAX'].values[0]), xy=(eda2_2_x['MAX'].valu
56 for spine in axs_right.spines.values():
57     spine.set_color('None')
58 plt.xlabel('\nPURCHASES_TRX', **xy_label)
59 plt.ylabel('')

```

```
60 plt.xticks(fontsize=8)
61 plt.yticks(fontsize=8)
62 plt.grid(axis='y', alpha=0)
63 plt.grid(axis='x', which='major', **grid_style)
64 plt.title('\nPurchase Total Transactions\n', **title)
65
66 plt.show();
```



EDA 3: Monto invertido en Compras a Cuotas vs Limite de credito del cliente

In [280]:

```

1 # --- EDA 3 Variables ---
2 scatter_style=dict(linewidth=0.65, edgecolor='#100C07', alpha=0.85)
3 sub_scatter_style_color=dict(s=5, alpha=0.65, linewidth=0.15, zorder=10, edgecolor='#
4 sub_scatter_style_grey=dict(s=5, alpha=0.3, linewidth=0.7, zorder=5, color='#CAC9CD')
5 grid_style=dict(alpha=0.3, color='#9B9A9C', linestyle='dotted', zorder=1)
6 xy_label=dict(fontweight='bold', fontsize=14, fontfamily='serif')
7 suptitle=dict(fontsize=22, fontweight='heavy', fontfamily='serif')
8 title=dict(fontsize=16, fontfamily='serif')
9 color_pallete=['#2D0F51', '#FF9A00', '#6600A5', '#FFD61E', '#722E9A', '#FFE863', '#A4
10 sub_axes=[None] * 7
11
12 # --- EDA 3 Data Frame ---
13 eda3 = df[['CREDIT_LIMIT', 'INSTALLMENTS_PURCHASES', 'TENURE']]
14 eda3['TENURE'] = eda1['TENURE'].astype(str)
15 tenure = sorted(eda1['TENURE'].unique())
16
17 # --- EDA 3 Settings ---
18 fig = plt.figure(figsize=(15, 20))
19 gs = fig.add_gridspec(7, 7)
20 ax = fig.add_subplot(gs[:7, :])
21 ax.set_aspect(1)
22
23 # --- EDA 3: Main Scatter Plot ---
24 for x in range(len(tenure)):
25     eda3_x = eda3[eda3['TENURE']==tenure[x]]
26     ax.scatter(eda3_x['CREDIT_LIMIT'], eda3_x['INSTALLMENTS_PURCHASES'], s=80, color=
27     ax.set_title('There is no heteroscedasticity detected between the credit limit an
28     ax.set_xlabel('\nCREDIT_LIMIT', **xy_label)
29     ax.set_ylabel('INSTALLMENTS_PURCHASES\n', **xy_label)
30     ax.grid(axis='y', which='major', **grid_style)
31     ax.grid(axis='x', which='major', **grid_style)
32     for spine in ax.spines.values():
33         spine.set_color('None')
34     for spine in ['bottom', 'left']:
35         ax.spines[spine].set_visible(True)
36         ax.spines[spine].set_color('#CAC9CD')
37     plt.xticks(fontsize=12)
38     plt.yticks(fontsize=12)
39
40 # --- EDA 3: Sub Plots ---
41 for idx, tnr in enumerate(tenure):
42     sub_axes[idx] = fig.add_subplot(gs[6, idx], aspect=1)
43
44     sub_axes[idx].scatter(eda3[eda3['TENURE']!=tnr]['CREDIT_LIMIT'], eda3[eda3['TENUR
45     sub_axes[idx].scatter(eda3[eda3['TENURE']==tnr]['CREDIT_LIMIT'], eda3[eda3['TENUR
46
47     cnt = (eda3['TENURE']==tnr).sum()
48     sub_axes[idx].set_title(f'Tenure {tnr} - ({cnt})', loc='left', fontsize=10, fontf
49     sub_axes[idx].set_xticks([])
50     sub_axes[idx].set_yticks([])
51     for spine in sub_axes[idx].spines.values():
52         spine.set_color('None')
53
54 # --- EDA 3 XY Limit ---
55 for axes in [ax] + sub_axes:
56     axes.set_xlim(-1000, 31000)
57     axes.set_ylim(-1000, 25000)
58
59 # --- EDA 3 Title ---

```

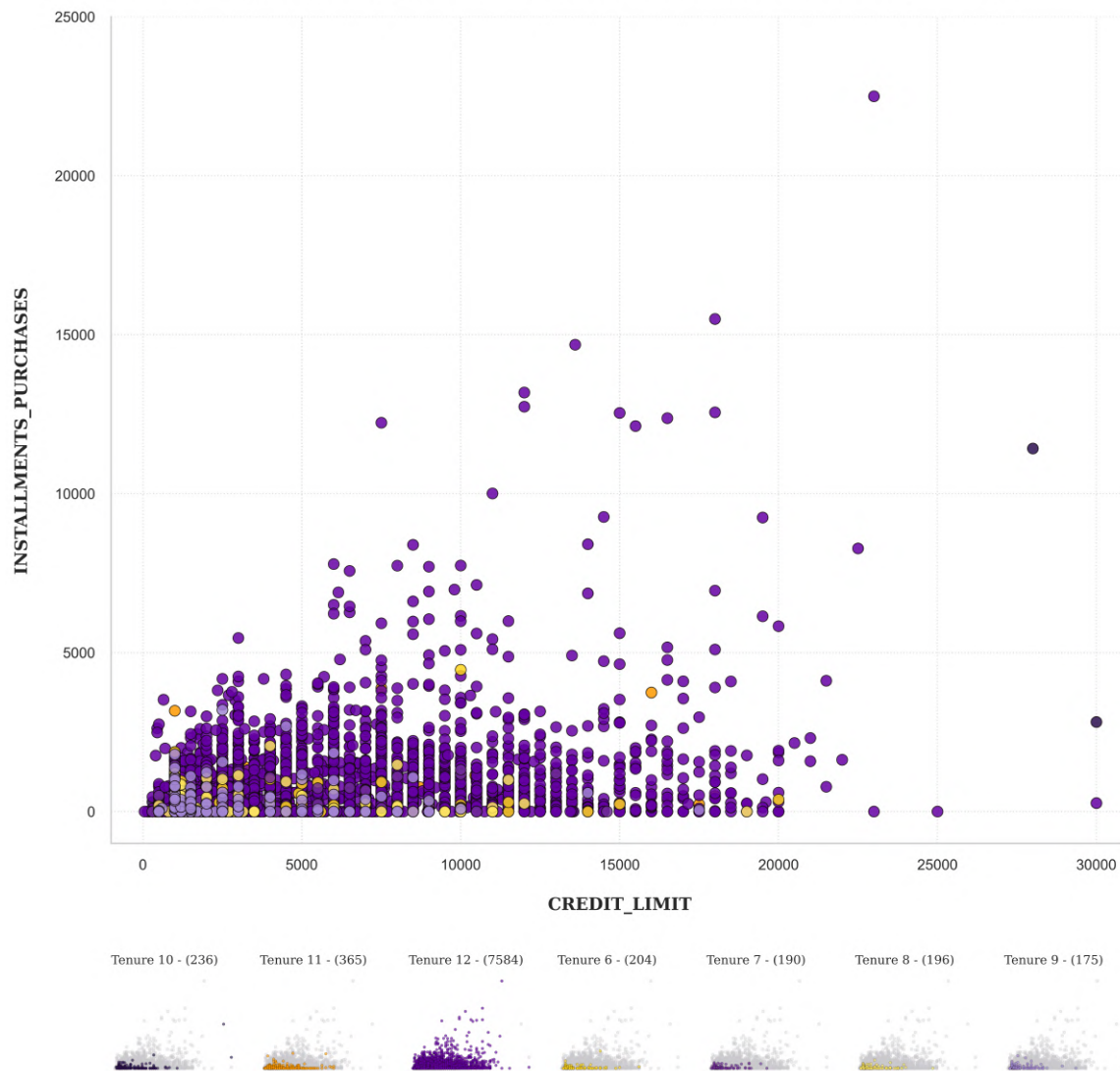
```

60 plt.suptitle('Credit Limit vs. Installment Purchases based on Tenure', x=0.123, y=0.7
61
62 plt.show();

```

Credit Limit vs. Installment Purchases based on Tenure

There is no heteroscedasticity detected between the credit limit and installment purchases.



Preprocesado de los datos

Solo debemos tener variables numericas hasta este punto

In [281]:

```

1 # --- Showing Dataframe ---
2 df = df #xd
3 df.head().style.background_gradient(cmap='afmhot_r')

```

Out[281]:

	BALANCE	BALANCE_FREQUENCY	PURCHASES	ONEOFF_PURCHASES	INSTALLMENTS
0	40.900749	0.818182	95.400000	0.000000	
1	3202.467416	0.909091	0.000000	0.000000	
2	2495.148862	1.000000	773.170000	773.170000	
3	1666.670542	0.636364	1499.000000	1499.000000	
4	817.714335	1.000000	16.000000	16.000000	

Imputation

Como el conjunto de datos es para clustering, se utilizará KNNImputer() para evitar resultados de clustering sesgados. El valor medio de los n_neighbors más cercanos encontrados en el conjunto de datos se utiliza para imputar los valores faltantes para cada muestra.

Las columnas 'credit_limit' y 'minimum_payments' contienen valores nulos que deben ser cargados con valores antes de continuar.

In [282]:

```

1 # --- List Null Columns ---
2 null_columns = df.columns[df.isnull().any()].tolist()
3
4 # --- Perform Imputation ---
5 imputer = KNNImputer()
6 df_imp = pd.DataFrame(imputer.fit_transform(df[null_columns]), columns=null_columns)
7 df = df.fillna(df_imp)
8
9 # --- Showing Dataframe ---
10 print clr.start + '.: Dataframe after Imputation :.' + clr.end
11 print clr.color + '*' * 33
12 df.head().style.background_gradient(cmap='afmhot_r')

```

```

.: Dataframe after Imputation :.
*****

```

Out[282]:

	BALANCE	BALANCE_FREQUENCY	PURCHASES	ONEOFF_PURCHASES	INSTALLMENTS
0	40.900749	0.818182	95.400000	0.000000	
1	3202.467416	0.909091	0.000000	0.000000	
2	2495.148862	1.000000	773.170000	773.170000	
3	1666.670542	0.636364	1499.000000	1499.000000	
4	817.714335	1.000000	16.000000	16.000000	

Scaling

El siguiente paso es escalar el conjunto de datos. La escalación es esencial ya que maneja la variabilidad del conjunto de datos, transforma los datos en un rango definido utilizando una transformación lineal para producir clusters de alta calidad y mejora la precisión de los algoritmos de clustering. En este caso, se utilizará un escalador estándar para estandarizar las características eliminando la media y escalando a la varianza unitaria.

In [283]:

```
1 # --- Scaling Dataset w/ Standard Scaler ---  
2 X = pd.DataFrame(StandardScaler().fit_transform(df))
```

Hopkins Test

El siguiente paso es realizar una prueba estadística utilizando la prueba estadística de Hopkins para el conjunto de datos preprocesado para medir la tendencia de clustering de los datos (medida de en qué grado existen clusters en los datos a ser agrupados).

A continuación se muestra la hipótesis de la prueba estadística de Hopkins: H0: El conjunto de datos no está distribuido uniformemente (contiene clusters significativos). H1: El conjunto de datos está distribuido uniformemente (no hay clusters significativos). Criterios: Si el valor está entre {0.7, ..., 0.99}, se acepta H0 (tiene una alta tendencia a agruparse).

In [284]:

```

1  # --- Hopkins Test (codes by Matevž Kunaver) ---
2  def hopkins(X):
3      d = X.shape[1]
4      n = len(X)
5      m = int(0.1 * n)
6      nbrs = NearestNeighbors(n_neighbors=1).fit(X)
7
8      rand_X = sample(range(0, n, 1), m)
9
10     ujd = []
11     wjd = []
12     for j in range(0, m):
13         u_dist, _ = nbrs.kneighbors(uniform(np.amin(X,axis=0),np.amax(X,axis=0),d).re
14         ujd.append(u_dist[0][1])
15         w_dist, _ = nbrs.kneighbors(X.iloc[rand_X[j]].values.reshape(1, -1), 2, retur
16         wjd.append(w_dist[0][1])
17
18     H = sum(ujd) / (sum(ujd) + sum(wjd))
19     if isnan(H):
20         print (ujd, wjd)
21         H = 0
22
23     return H
24
25 # --- Perform Hopkins Test ---
26 hopkins_value = hopkins(X)
27 hopkins_result = 'Result: '+clr.start+'{:.4f}'.format(hopkins_value)+clr.end
28 print(clr.start+'.: Hopkins Test :'+clr.end)
29 print(clr.color+'*' * 19+clr.end)
30 print(hopkins_result)
31 if 0.7 < hopkins_value < 0.99:
32     print('>> From the result above,'+clr.color+' it has a high tendency to cluster (
33     print('\n'+clr.color+'*' * 31+clr.end)
34     print(clr.start+'... Conclusions: Accept H0 ...'+clr.end)
35     print(clr.color+'*' * 31+clr.end)
36 else:
37     print('>> From the result above,'+clr.color+' it has no meaningful clusters'+clr.
38     print('\n'+clr.color+'*' * 31+clr.end)
39     print(clr.start+'... Conclusions: Reject H0 ...'+clr.end)
40     print(clr.color+'*' * 31+clr.end)

```

.: Hopkins Test .:

Result: 0.9666

>> From the result above, it has a high tendency to cluster (contains mean
ingful clusters)

... Conclusions: Accept H0 ...

PCA

Debido a que el conjunto de datos tiene 18 columnas, es fundamental aplicar algún método de reducción de dimensionalidad. Sin duda, uno de los métodos más utilizados es el PCA. El análisis de componentes principales (PCA) es un método utilizado en el aprendizaje automático no supervisado (como el

agrupamiento) que reduce los datos de alta dimensión a dimensiones más pequeñas mientras se preserva tanta información como sea posible.

- En este notebook, el número de características se reducirá a 2 dimensiones para que los resultados del clustering puedan ser visualizados.

In [285]:

```
1 # --- Transform into Array ---
2 X = np.asarray(X)
3
4 # --- Applying PCA ---
5 pca = PCA(n_components=2, random_state=24)
6 X = pca.fit_transform(X)
```

Clusterización de clientes

En esta sección, implementaremos las técnicas de clustering mencionadas en la primera sección. Además, se proporcionará una explicación para cada modelo.

7.1 | K-Means

Este algoritmo de aprendizaje no supervisado agrupa datos en clusters, con el objetivo de minimizar la variación dentro de cada cluster. Primero se elige el número de clusters (k) que se desea crear, y luego se seleccionan k puntos iniciales aleatorios como centroides. Luego, cada punto de datos se asigna al cluster cuyo centroide está más cerca. Los centroides se recalculan en función de los puntos que se han asignado a su cluster, y el proceso se repite hasta que no hay cambios en la asignación de los puntos de datos.

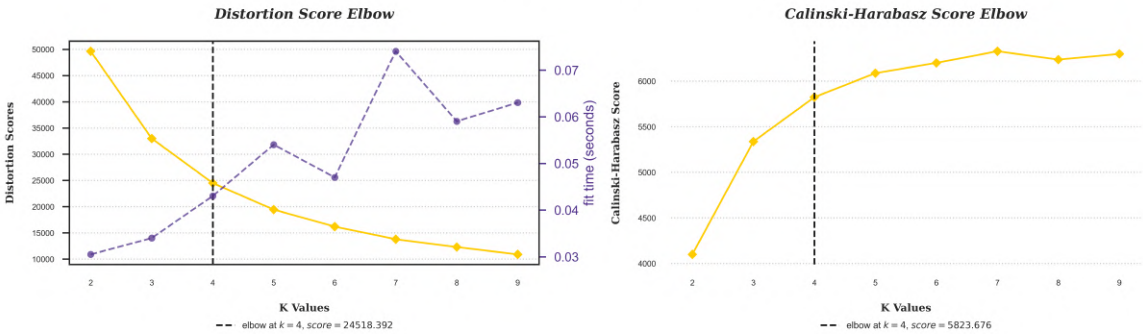
In [286]:

```

1  # --- Define K-Means Functions ---
2  def kmeans():
3
4      # --- Figures Settings ---
5      color_palette=['#FFCC00', '#54318C']
6      set_palette(color_palette)
7      title=dict(fontsize=12, fontweight='bold', style='italic', fontfamily='serif')
8      text_style=dict(fontweight='bold', fontfamily='serif')
9      fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(14, 5))
10
11     # --- Elbow Score ---
12     elbow_score = KElbowVisualizer(KMeans(random_state=32, max_iter=500), k=(2, 10),
13     elbow_score.fit(X)
14     elbow_score.finalize()
15     elbow_score.ax.set_title('Distortion Score Elbow\n', **title)
16     elbow_score.ax.tick_params(labelsize=7)
17     for text in elbow_score.ax.legend_.texts:
18         text.set_fontsize(9)
19     for spine in elbow_score.ax.spines.values():
20         spine.set_color('None')
21     elbow_score.ax.legend(loc='upper center', bbox_to_anchor=(0.5, -0.15), borderpad=
22     elbow_score.ax.grid(axis='y', alpha=0.5, color='#9B9A9C', linestyle='dotted')
23     elbow_score.ax.grid(axis='x', alpha=0)
24     elbow_score.ax.set_xlabel('\nK Values', fontsize=9, **text_style)
25     elbow_score.ax.set_ylabel('Distortion Scores\n', fontsize=9, **text_style)
26
27     # --- Elbow Score (Calinski-Harabasz Index) ---
28     elbow_score_ch = KElbowVisualizer(KMeans(random_state=32, max_iter=500), k=(2, 10
29     elbow_score_ch.fit(X)
30     elbow_score_ch.finalize()
31     elbow_score_ch.ax.set_title('Calinski-Harabasz Score Elbow\n', **title)
32     elbow_score_ch.ax.tick_params(labelsize=7)
33     for text in elbow_score_ch.ax.legend_.texts:
34         text.set_fontsize(9)
35     for spine in elbow_score_ch.ax.spines.values():
36         spine.set_color('None')
37     elbow_score_ch.ax.legend(loc='upper center', bbox_to_anchor=(0.5, -0.15), borderp
38     elbow_score_ch.ax.grid(axis='y', alpha=0.5, color='#9B9A9C', linestyle='dotted')
39     elbow_score_ch.ax.grid(axis='x', alpha=0)
40     elbow_score_ch.ax.set_xlabel('\nK Values', fontsize=9, **text_style)
41     elbow_score_ch.ax.set_ylabel('Calinski-Harabasz Score\n', fontsize=9, **text_styl
42
43     plt.suptitle('Credit Card Customer Clustering using K-Means', fontsize=14, **text
44
45     plt.tight_layout()
46     plt.show();
47
48     # --- Calling K-Means Functions ---
49     kmeans();

```

Credit Card Customer Clustering using K-Means



In [259]:

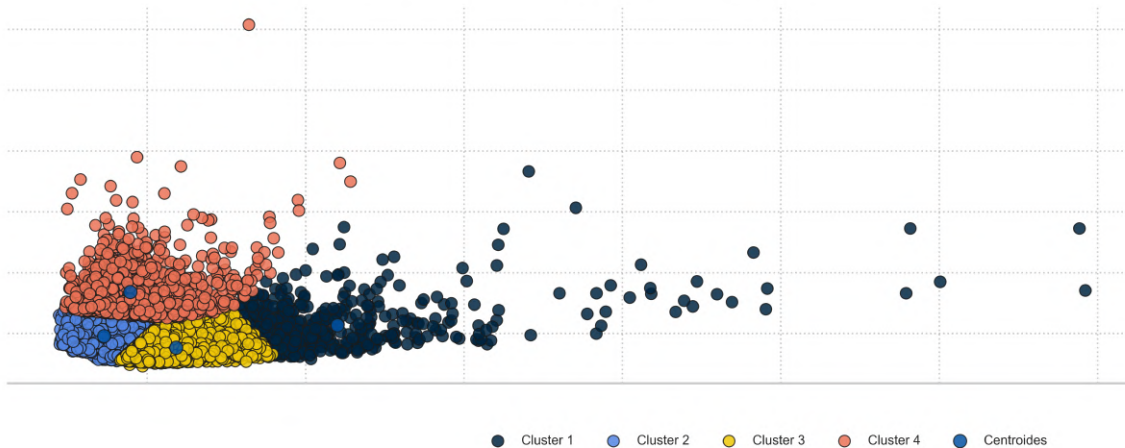
```

1 # --- Implementing K-Means ---
2 kmeans = KMeans(n_clusters=4, random_state=32, max_iter=500)
3 y_kmeans = kmeans.fit_predict(X)
4
5 # --- Define K-Means Visualizer & Plots ---
6 def visualizer(kmeans, y_kmeans):
7
8     # --- Figures Settings ---
9     #cluster_colors=['#FFBB00', '#3C096C', '#9D4EDD', '#FFE270']
10    cluster_colors = ['#002642', '#5386E4', '#F0C808', '#EC7357']
11
12
13
14    labels = ['Cluster 1', 'Cluster 2', 'Cluster 3', 'Cluster 4', 'Centroides']
15    title=dict(fontsize=12, fontweight='bold', style='italic', fontfamily='serif')
16    text_style=dict(fontweight='bold', fontfamily='serif')
17    scatter_style=dict(linewidth=0.65, edgecolor='#1B1B1B', alpha=0.85)
18    legend_style=dict(borderpad=2, frameon=False, fontsize=8)
19    fig, (ax2, ax3) = plt.subplots(2, 1, figsize=(10, 10))
20
21
22    y_kmeans_labels = list(set(y_kmeans.tolist()))
23    for i in y_kmeans_labels:
24        ax2.scatter(X[y_kmeans==i, 0], X[y_kmeans == i, 1], s=50, c=cluster_colors[i])
25    ax2.scatter(kmeans.cluster_centers_[0], kmeans.cluster_centers_[1], s=65, c=cluster_colors[0])
26    for spine in ax2.spines.values():
27        spine.set_color('None')
28    ax2.set_title('Scatter Plot - Distribucion en Clusters\n', **title)
29    ax2.legend(labels, bbox_to_anchor=(0.95, -0.05), ncol=5, **legend_style)
30    ax2.grid(axis='both', alpha=0.5, color='#9B9A9C', linestyle='dotted')
31    ax2.tick_params(left=False, right=False, labelleft=False, labelbottom=False, bottom=False)
32    ax2.spines['bottom'].set_visible(True)
33    ax2.spines['bottom'].set_color('#CAC9CD')
34
35    # --- Waffle Chart ---
36    unique, counts = np.unique(y_kmeans, return_counts=True)
37    df_waffle = dict(zip(unique, counts))
38    total = sum(df_waffle.values())
39    wfl_square = {key: value/100 for key, value in df_waffle.items()}
40    wfl_label = {key: round(value/total*100, 2) for key, value in df_waffle.items()}
41
42    ax3=plt.subplot(2, 2, (3,4))
43    ax3.set_title('Porcentaje de Clientes en cada Cluster\n', **title)
44    ax3.set_aspect(aspect='auto')
45    Waffle.make_waffle(ax=ax3, rows=8, values=wfl_square, colors=cluster_colors,
46                        labels=[f"Cluster {i+1} - ({k}%) " for i, k in wfl_label.items()],
47                        legend={'loc': 'upper center', 'bbox_to_anchor': (0.5, -0.05),
48                               'frameon': False, 'fontsize':10})
49    ax3.text(0.01, -0.09, '** 1 cuadro ≈ 100 Clientes', weight = 'bold', style='italic')
50
51    # --- Suptitle & WM ---
52    plt.suptitle('Clientes de Tarjeta de Credito - Clustering via K-Means\n', fontsize=12)
53
54    plt.tight_layout()
55    plt.show();
56
57 # --- Calling K-Means Functions ---

```

Clientes de Tarjeta de Credito - Clustering via K-Means

Scatter Plot - Distribucion en Clusters



Porcentaje de Clientes en cada Cluster



In [287]:

```
1 # --- Evaluate Clustering Quality Function ---
2 def evaluate_clustering(X, y):
3     db_index = round(davies_bouldin_score(X, y), 3)
4     s_score = round(silhouette_score(X, y), 3)
5     ch_index = round(calinski_harabasz_score(X, y), 3)
6     print clr.start + '.: Evaluate Clustering Quality :.' + clr.end
7     print clr.color + '*' * 34 + clr.end
8     print('.: Davies-Bouldin Index: ' + clr.start, db_index)
9     print(clr.end + '.: Silhouette Score: ' + clr.start, s_score)
10    print(clr.end + '.: Calinski Harabasz Index: ' + clr.start, ch_index)
11    return db_index, s_score, ch_index
12
13 # --- Evaluate K-Means Cluster Quality ---
14 db_kmeans, ss_kmeans, ch_kmeans = evaluate_clustering(X, y_kmeans)
```

.: Evaluate Clustering Quality :.

.: Davies-Bouldin Index: 0.801

.: Silhouette Score: 0.408

.: Calinski Harabasz Index: 5823.676

- 📌 El índice Davies-Bouldin es una métrica para evaluar algoritmos de agrupamiento. Se define como una relación entre la dispersión del cluster y la separación del cluster. Los puntajes oscilan entre 0 y

más. 0 indica un mejor agrupamiento.

- 📌 El Coeficiente/Puntuación de Silueta es una métrica utilizada para calcular la bondad de una técnica de agrupamiento. Su valor varía de -1 a 1. Cuanto mayor sea la puntuación, mejor. 1 significa que los clusters están bien separados entre sí y claramente distinguibles. Cero significa que los clusters son indiferentes/la distancia entre clusters no es significativa. -1 significa que los clusters están asignados de manera incorrecta.
- 📌 El índice Calinski-Harabasz (también conocido como el criterio de relación de varianza) es la relación entre la suma de la dispersión entre clusters y la dispersión dentro de los clusters para todos los clusters. Cuanto mayor sea la puntuación, mejores serán los resultados.

7.2 | DBSCAN

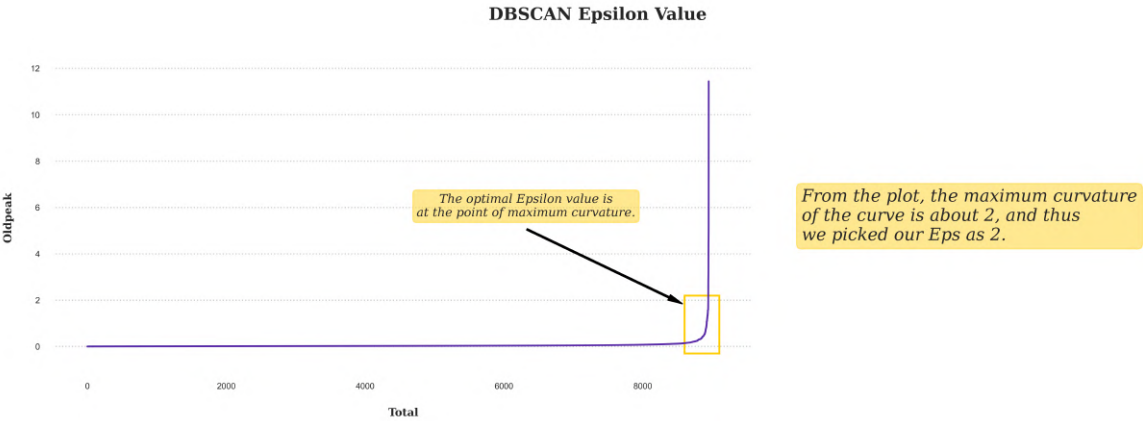
Este algoritmo también es de aprendizaje no supervisado y se utiliza para agrupar puntos de datos en función de su densidad. El algoritmo comienza con un punto aleatorio y luego busca todos los puntos cercanos que estén dentro de una distancia específica (epsilon). Luego, se identifican los grupos de puntos que están dentro de un radio específico (MinPoints). Los puntos que no están dentro de ningún grupo se consideran valores atípicos.

In [288]:

```

1 # --- Define Epsilon Values ---
2 def epsilon():
3
4     # --- Calculate Nearest Neighbors ---
5     neighbors=NearestNeighbors(n_neighbors=2)
6     nbrs=neighbors.fit(X)
7     distances, indices=nbrs.kneighbors(X)
8     distances=np.sort(distances, axis = 0)
9
10    # --- Figure Settings ---
11    bbox=dict(boxstyle='round', pad=0.3, color='#FFDA47', alpha=0.6)
12    txt1=dict(textcoords='offset points', va='center', ha='center', fontfamily='serif')
13    txt2=dict(textcoords='offset points', va='center', fontfamily='serif', style='italic')
14    kw=dict(arrowstyle='Simple, tail_width=0.1, head_width=0.4, head_length=1', color='red')
15    text_style=dict(fontweight='bold', fontfamily='serif')
16    fig=plt.figure(figsize=(14, 5))
17
18    # --- Epsilon Plot ---
19    distances_1=distances[:, 1]
20    ax1=fig.add_subplot(1, 3, (1, 2))
21    plt.plot(distances_1, color='#5829A7')
22    plt.xlabel('\nTotal', fontsize=9, **text_style)
23    plt.ylabel('Oldpeak\n', fontsize=9, **text_style)
24    ax1.add_patch(Rectangle((8600, -0.3), 500, 2.5, edgecolor='#FFCC00', fill=False,
25    plt.annotate('The optimal Epsilon value is\nat the point of maximum curvature.',
26    plt.annotate('', xy=(8600, 1.8), xytext=(6300, 5.1), arrowprops=kw)
27    for spine in ax1.spines.values():
28        spine.set_color('None')
29    plt.grid(axis='y', alpha=0.5, color='#9B9A9C', linestyle='dotted')
30    plt.grid(axis='x', alpha=0)
31    plt.tick_params(labelsize=7)
32
33    # --- Explanations ---
34    ax2=fig.add_subplot(1, 3, 3)
35    plt.annotate('From the plot, the maximum curvature\nof the curve is about 2, and
36    for spine in ax2.spines.values():
37        spine.set_color('None')
38    plt.grid(axis='both', alpha=0)
39    plt.axis('off')
40
41    plt.suptitle('DBSCAN Epsilon Value\n', fontsize=14, **text_style)
42
43    plt.tight_layout()
44    plt.show();
45
46 # --- Calling Epsilon Functions ---
47 epsilon();

```



In [299]:

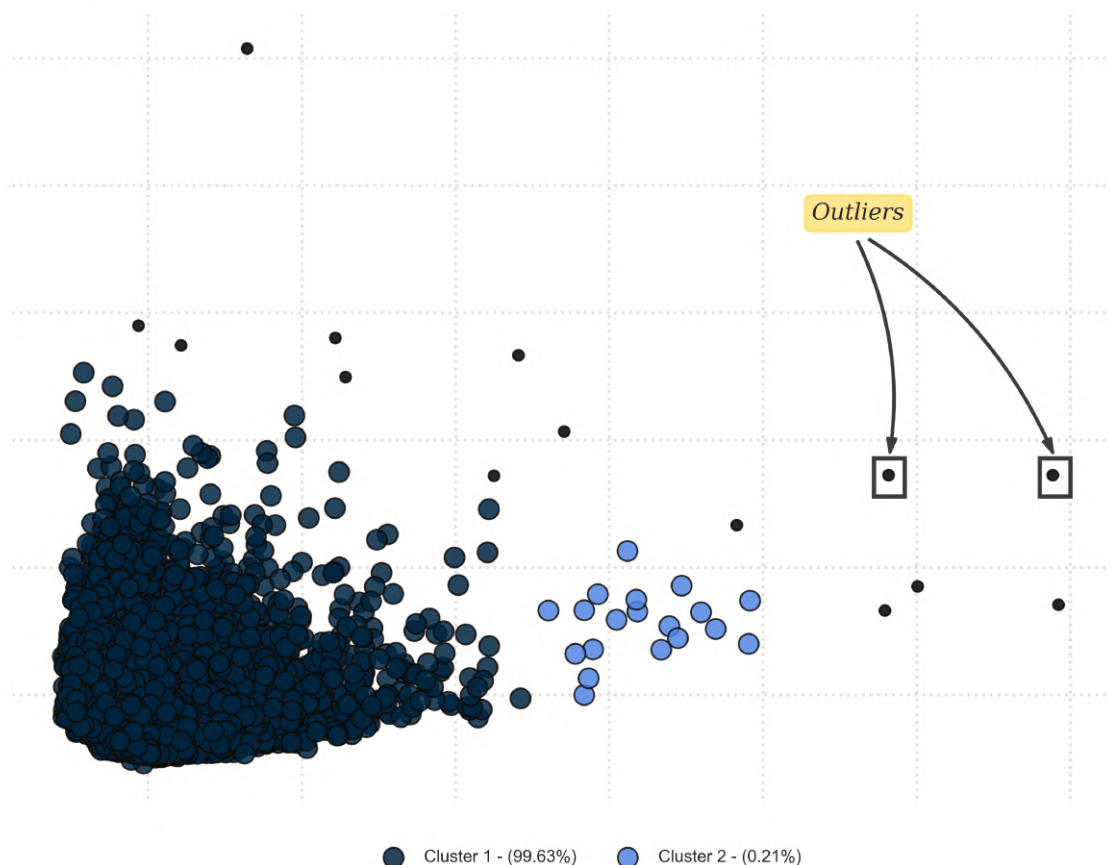
```

1 # --- Implementing DBSCAN ---
2 dbscan = DBSCAN(eps=2, min_samples=4)
3 y_dbscan = dbscan.fit_predict(X)
4
5 # --- Define DBSCAN Result Distribution ---
6 def dbscan_visualizer(dbscan, y_dbscan):
7
8     # --- Figures Settings ---
9     #cluster_colors=['#FFBB00', '#9D4EDD', 'black']
10    cluster_colors = ['#002642', '#5386E4', 'black']
11    labels = ['Cluster 1', 'Cluster 2', 'Outliers']
12    suptitle=dict(fontsize=12, fontweight='heavy', fontfamily='serif')
13    title=dict(fontsize=8, fontfamily='serif')
14    scatter_style=dict(linewidth=0.65, edgecolor='#100C07', alpha=0.85)
15    bbox=dict(boxstyle='round', pad=0.3, color='#FFDA47', alpha=0.6)
16    txt=dict(textcoords='offset points', va='center', ha='center', fontfamily='serif')
17    legend_style=dict(borderpad=2, frameon=False, fontsize=6)
18
19    # --- Arrow Settings ---
20    style = 'Simple, tail_width=0.3, head_width=3, head_length=5'
21    kw = dict(arrowstyle=style, color='#3E3B39')
22    arrow1 = patches.FancyArrowPatch((23, 18), (24.1, 9.3), connectionstyle='arc3, ra
23    arrow2 = patches.FancyArrowPatch((23.3, 18), (29.5, 9.3), connectionstyle='arc3,
24
25    # --- Percentage Labels ---
26    unique, counts = np.unique(y_dbscan, return_counts=True)
27    dbscan_count = dict(zip(unique, counts))
28    total = sum(dbscan_count.values())
29    dbscan_label = {key: round(value/total*100, 2) for key, value in dbscan_count.items()}
30
31    # --- Clusters Distribution ---
32    y_dbscan_labels = list(set(y_dbscan.tolist()))
33    fig, ax = plt.subplots(1, 1, figsize=(7, 5))
34    for i in np.arange(0, 2, 1):
35        plt.scatter(X[y_dbscan==i, 0], X[y_dbscan == i, 1], s=50, c=cluster_colors[i])
36    plt.scatter(X[y_dbscan==-1, 0], X[y_dbscan == -1, 1], s=15, c=cluster_colors[2],
37    for spine in ax.spines.values():
38        spine.set_color('None')
39    plt.legend([f"Cluster {i+1} - ({k}%)" for i, k in dbscan_label.items()], bbox_to_
40    plt.grid(axis='both', alpha=0.3, color='#9B9A9C', linestyle='dotted')
41    ax.add_patch(Rectangle((29, 7.8), 1, 1.5, edgecolor='#3E3B39', fill=False, lw=1.5
42    ax.add_patch(Rectangle((23.6, 7.8), 1, 1.5, edgecolor='#3E3B39', fill=False, lw=1
43    ax.add_patch(arrow1)
44    ax.add_patch(arrow2)
45    plt.annotate('Outliers', xy=(23, 18.8), xytext=(1, 1), fontsize=8, bbox=bbox, **t
46    plt.tick_params(left=False, right=False, labelleft=False, labelbottom=False, bc
47    plt.title('Two clusters of credit card customers were formed. There are also some
48    plt.suptitle('Credit Card Customer Clustering using DBSCAN', x=0.123, y=0.98, ha=
49
50    plt.show();
51
52 # --- Calling DBSCAN Functions ---
53 dbscan_visualizer(dbscan, y_dbscan);

```

Credit Card Customer Clustering using DBSCAN

Two clusters of credit card customers were formed. There are also some outliers detected.



In [193]:

```
1 # --- Evaluate DBSCAN Cluster Quality ---
2 db_dbscan, ss_dbscan, ch_dbscan = evaluate_clustering(X, y_dbscan)
```

.: Evaluate Clustering Quality :.

.: Davies-Bouldin Index: **1.287**

.: Silhouette Score: **0.803**

.: Calinski Harabasz Index: **685.303**

7.3 | Hierarchical Clustering (Agglomerative)

Hierarchical Clustering (Agglomerative): Este algoritmo también agrupa datos en clusters, pero utiliza una técnica jerárquica. En primer lugar, cada punto de datos se considera un cluster individual. Luego, se unen los dos clusters más cercanos, y este proceso se repite hasta que todos los puntos de datos se agrupan en un solo cluster. Se puede visualizar el proceso en un dendrograma, que muestra la estructura jerárquica de los clusters.

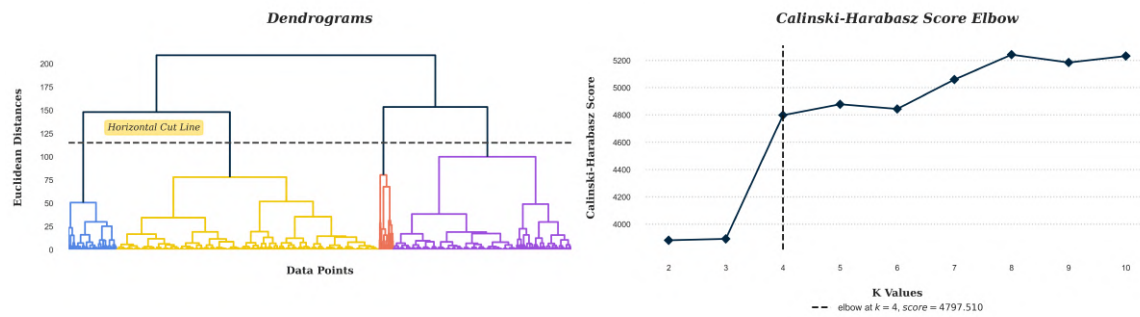
In [300]:

```

1 # --- Define Dendrogram ---
2 def agg_dendrogram():
3
4     # --- Figure Settings ---
5     #color_palette=['#472165', '#FFBB00', '#3C096C', '#9D4EDD', '#FFE270']
6     color_palette = ['#002642', '#5386E4', '#F0C808', '#EC7357', '#9D4EDD']
7     set_palette(color_palette)
8     text_style=dict(fontweight='bold', fontfamily='serif')
9     ann=dict(textcoords='offset points', va='center', ha='center', fontfamily='serif')
10    title=dict(fontsize=12, fontweight='bold', style='italic', fontfamily='serif')
11    bbox=dict(boxstyle='round', pad=0.3, color='#FFDA47', alpha=0.6)
12    fig=plt.figure(figsize=(14, 5))
13
14    # --- Dendrogram Plot ---
15    ax1=fig.add_subplot(1, 2, 1)
16    dend=shc.dendrogram(shc.linkage(X, method='ward', metric='euclidean'))
17    plt.axhline(y=115, color='#3E3B39', linestyle='--')
18    plt.xlabel('\nData Points', fontsize=9, **text_style)
19    plt.ylabel('Euclidean Distances\n', fontsize=9, **text_style)
20    plt.annotate('Horizontal Cut Line', xy=(15000, 130), xytext=(1, 1), fontsize=8, b
21    plt.tick_params(labelbottom=False)
22    for spine in ax1.spines.values():
23        spine.set_color('None')
24    plt.grid(axis='both', alpha=0)
25    plt.tick_params(labelsize=7)
26    plt.title('Dendrograms\n', **title)
27
28    # --- Elbow Score (Calinski-Harabasz Index) ---
29    ax2=fig.add_subplot(1, 2, 2)
30    elbow_score_ch = KElbowVisualizer(AgglomerativeClustering(), metric='calinski_har
31    elbow_score_ch.fit(X)
32    elbow_score_ch.finalize()
33    elbow_score_ch.ax.set_title('Calinski-Harabasz Score Elbow\n', **title)
34    elbow_score_ch.ax.tick_params(labelsize=7)
35    for text in elbow_score_ch.ax.legend_.texts:
36        text.set_fontsize(9)
37    for spine in elbow_score_ch.ax.spines.values():
38        spine.set_color('None')
39    elbow_score_ch.ax.legend(loc='upper center', bbox_to_anchor=(0.5, -0.15), borderp
40    elbow_score_ch.ax.grid(axis='y', alpha=0.5, color='#9B9A9C', linestyle='dotted')
41    elbow_score_ch.ax.grid(axis='x', alpha=0)
42    elbow_score_ch.ax.set_xlabel('\nK Values', fontsize=9, **text_style)
43    elbow_score_ch.ax.set_ylabel('Calinski-Harabasz Score\n', fontsize=9, **text_styl
44
45    plt.suptitle('Credit Card Customer Clustering using Hierarchical Clustering\n', f
46
47    plt.tight_layout()
48    plt.show();
49
50 # --- Calling Dendrogram Functions ---
51 agg_dendrogram();

```

Credit Card Customer Clustering using Hierarchical Clustering



Basándonos en la distancia euclidiana en el dendrograma de arriba, se puede concluir que el número de clusters será cuatro ya que la línea vertical más alta/mayor distancia se encuentra en la primera línea/rama (a la izquierda de la imagen) y el umbral corta el dendrograma en cuatro partes. Además, según la puntuación de Calinski-Harabasz, el número óptimo de clusters obtenido es 4.

A continuación, implementaremos este número en el algoritmo de clustering aglomerativo y visualizaremos y evaluaremos los clusters creados.

In [301]:

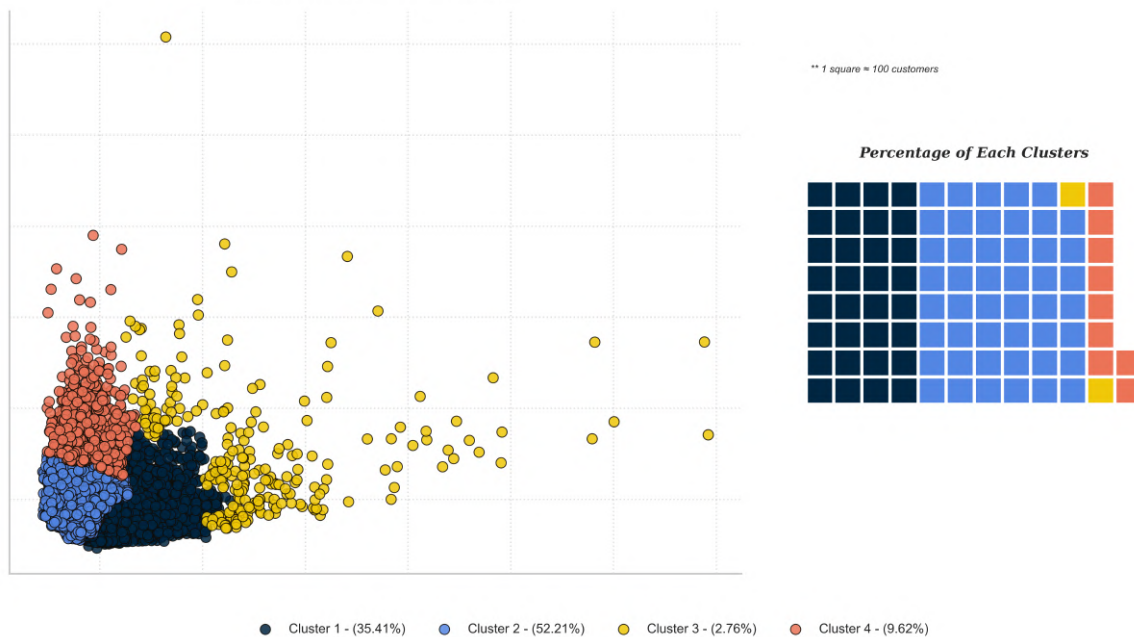
```

1 # --- Implementing Hierarchical Clustering ---
2 agg_cluster = AgglomerativeClustering(n_clusters=4, affinity='euclidean', linkage='wa
3 y_agg_cluster = agg_cluster.fit_predict(X)
4
5 # --- Define Hierarchical Clustering Distributions ---
6 def agg_visualizer(agg_cluster, y_agg_cluster):
7
8     # --- Figures Settings ---
9     cluster_colors=['#002642', '#5386E4', '#F0C808', '#EC7357']
10
11
12     labels = ['Cluster 1', 'Cluster 2', 'Cluster 3', 'Cluster 4']
13     suptitle=dict(fontsize=14, fontweight='heavy', fontfamily='serif')
14     title=dict(fontsize=10, fontweight='bold', style='italic', fontfamily='serif')
15     scatter_style=dict(linewidth=0.65, edgecolor='#100C07', alpha=0.85)
16     legend_style=dict(borderpad=2, frameon=False, fontsize=9)
17     fig=plt.figure(figsize=(14, 7))
18
19     # --- Percentage Labels ---
20     unique, counts = np.unique(y_agg_cluster, return_counts=True)
21     df_waffle = dict(zip(unique, counts))
22     total = sum(df_waffle.values())
23     wfl_square = {key: value/100 for key, value in df_waffle.items()}
24     wfl_label = {key: round(value/total*100, 2) for key, value in df_waffle.items()}
25
26     # --- Clusters Distribution ---
27     y_agg_labels = list(set(y_agg_cluster.tolist()))
28     ax1=fig.add_subplot(1, 3, (1, 2))
29     for i in y_agg_labels:
30         ax1.scatter(X[y_agg_cluster==i, 0], X[y_agg_cluster == i, 1], s=50, c=cluster
31     for spine in ax1.spines.values():
32         spine.set_color('None')
33     for spine in ['bottom', 'left']:
34         ax1.spines[spine].set_visible(True)
35         ax1.spines[spine].set_color('#CAC9CD')
36     ax1.legend([f"Cluster {i+1} - ({k}%)" for i, k in wfl_label.items()], bbox_to_anc
37     ax1.grid(axis='both', alpha=0.3, color='#9B9A9C', linestyle='dotted')
38     ax1.tick_params(left=False, right=False, labelleft=False, labelbottom=False, bc
39     plt.title('Scatter Plot Clusters Distributions\n', **title)
40
41     # --- Waffle Chart ---
42     ax2=fig.add_subplot(1, 3, 3)
43     ax2.set_title('Percentage of Each Clusters\n', **title)
44     ax2.set_aspect(aspect='equal')
45     Waffle.make_waffle(ax=ax2, rows=8, values=wfl_square, colors=cluster_colors)
46     ax2.get_legend().remove()
47
48     ax2.text(0.01, 1.5, '** 1 square ≈ 100 customers', style='italic', fontsize=7)
49
50     plt.suptitle('Credit Card Customer Clustering using Hierarchical Clustering\n', *
51     plt.show();
52
53
54 # --- Calling Hierarchical Clustering Functions ---
55 agg_visualizer(agg_cluster, y_agg_cluster);

```


Credit Card Customer Clustering using Hierarchical Clustering

Scatter Plot Clusters Distributions



A partir de la implementación de clustering jerárquico, se puede ver que se formaron 4 grupos. De los 4 grupos, el grupo 2 tiene la mayoría de los puntos de datos, seguido por el grupo 1. Sin embargo, cuando se comparan con los resultados del clustering por K-Means, los resultados del grupo 2 utilizando jerárquico tienen un porcentaje más significativo. Además, el algoritmo de clustering jerárquico considera que los valores atípicos son parte del grupo 3.

El último paso es evaluar la calidad del clustering que ofrece el clustering jerárquico. Se utilizarán el coeficiente de silueta y el índice Davies-Bouldin para evaluar la calidad.

In [292]:

```
1 # --- Evaluate DBSCAN Cluster Quality ---
2 db_agg, ss_agg, ch_agg = evaluate_clustering(X, y_agg_cluster)
```

.: Evaluate Clustering Quality :.

.: Davies-Bouldin Index: 0.863

.: Silhouette Score: 0.388

.: Calinski Harabasz Index: 4797.51

Evaluando los modelos utilizados

El siguiente paso es evaluar la calidad del agrupamiento proporcionado por K-Means. La evaluación de calidad utilizará el índice Davies-Bouldin, la puntuación de silueta y el índice Calinski-Harabasz.

In [293]:

```

1 # --- Comparison Table ---
2 compare = pd.DataFrame({'Model': ['K-Means', 'DBSCAN', 'Hierarchical Clustering'],
3                             'Davies-Bouldin Index': [db_kmeans, db_dbscan, db_agg],
4                             'Silhouette Score': [ss_kmeans, ss_dbscan, ss_agg],
5                             'Calinski-Harabasz Index': [ch_kmeans, ch_dbscan, ch_agg]})
6
7 # --- Create Accuracy Comparison Table ---
8 print(clear.start+': Model Accuracy Comparison :'+clear.end)
9 print(clear.color+'*' * 32+clear.end)
10 compare.sort_values(by='Model', ascending=False).style.background_gradient(cmap='inferno')

```

```

.: Model Accuracy Comparison :.
*****

```

Out[293]:

	Model	Davies-Bouldin Index	Silhouette Score	Calinski-Harabasz Index
0	K-Means	0.801000	0.408000	5823.676000
2	Hierarchical Clustering	0.863000	0.388000	4797.510000
1	DBSCAN	1.287000	0.803000	685.303000

- La tabla anterior muestra que el algoritmo K-Means tiene el índice Davies-Bouldin más bajo en comparación con los otros dos algoritmos, por lo que se puede concluir que **K-Means tiene una calidad de agrupamiento decente en comparación con los otros dos algoritmos**.
- Sin embargo, según la puntuación de silueta, **K-Means tiene** la segunda puntuación de silueta más alta (hay algunos **clusters superpuestos formados con este algoritmo**).
- A partir de los resultados del índice Calinski-Harabasz, K-Means tiene el índice más alto en comparación con otros algoritmos. Esto indica que **K-Means funciona mejor y es más denso que otros algoritmos**.

Perfilado de clientes

8. | K-means results

Debido al performance de k-means, de aquí en adelante solo será utilizada su clusterización para el análisis y perfilado de los clientes.

En la siguiente gráfica, podemos ver cómo se agrupan los mismos teniendo en cuenta el PCA, también podemos ver cuál es el % de clientes en cada grupo.

In [294]:

```

1 # --- Add K-Means Prediction to Data Frame ----
2 df['cluster_result'] = y_kmeans+1
3 df['cluster_result'] = 'Cluster '+df['cluster_result'].astype(str)
4
5 # --- Calculationg Overall Mean from Current Data Frame ---
6 df_profile_overall = pd.DataFrame()
7
8 df_profile_overall['Overall'] = df.describe().loc[['mean']].T
9
10 # --- Summarize Mean of Each Clusters ---
11 df_cluster_summary = df.groupby('cluster_result').describe().T.reset_index().rename(columns={'cluster_result': 'Cluster'})
12 df_cluster_summary = df_cluster_summary[df_cluster_summary['Metrics'] == 'mean'].set_index('Cluster')
13
14
15 # --- Combining Both Data Frame ---
16 print(clr.start+'.: Summarize of Each Clusters :.'+clr.end)
17 print(cclr.color+'*' * 33)
18 df_profile = df_cluster_summary.join(df_profile_overall).reset_index()
19
20
21 df_profile.style.background_gradient(axis=1,cmap='YlOrBr')

```

.: Summarize of Each Clusters :.

Out[294]:

	Column Name	Metrics	Cluster 1	Cluster 2	Cluster 3
0	BALANCE	mean	3401.840056	1012.915503	824.180354
1	BALANCE_FREQUENCY	mean	0.988427	0.799792	0.919997
2	PURCHASES	mean	6894.613917	223.081279	1236.499406
3	ONEOFF_PURCHASES	mean	4511.889901	157.576608	621.738764
4	INSTALLMENTS_PURCHASES	mean	2383.916859	65.833785	614.965725
5	CASH_ADVANCE	mean	773.154467	614.588758	147.442197
6	PURCHASES_FREQUENCY	mean	0.954443	0.190099	0.862440
7	ONEOFF_PURCHASES_FREQUENCY	mean	0.726667	0.074901	0.301885
8	PURCHASES_INSTALLMENTS_FREQUENCY	mean	0.808946	0.109962	0.675116
9	CASH_ADVANCE_FREQUENCY	mean	0.084806	0.122459	0.030698
10	CASH_ADVANCE_TRX	mean	2.363817	2.264421	0.569732
11	PURCHASES_TRX	mean	82.902584	2.960949	21.317296
12	CREDIT_LIMIT	mean	9541.650099	3109.010550	4250.051845
13	PAYMENTS	mean	6723.271522	856.400147	1328.949020
14	MINIMUM_PAYMENTS	mean	1830.297811	589.839851	600.410547
15	PRC_FULL_PAYMENT	mean	0.288014	0.065965	0.282867
16	TENURE	mean	11.960239	11.364216	11.661695

Basándonos en la tabla anterior, se puede concluir que cada cluster tiene las siguientes características:

- **Cluster 1 (Full Player Customer):** Los clientes en este cluster son usuarios activos de la tarjeta de crédito del banco. Esto se puede observar por la frecuencia con la que el saldo cambia y el monto de los saldos es lo suficientemente alto en comparación con los otros clusters. Además, en comparación con los otros clusters, este cluster tiene valores promedio más altos en varios aspectos. Los clientes de tarjetas de crédito en este cluster también usan activamente las tarjetas de crédito para facilitar transacciones y pagos a plazos. Los adelantos en efectivo, las transacciones y los pagos a plazos en este cluster también ocurren con más frecuencia. La antigüedad relativamente alta también muestra que el puntaje crediticio en este cluster es muy bueno.
- **Cluster 2 (Usuarios Principiantes/Estudiantes):** En contraste con el cluster 1, los clientes rara vez o casi nunca usan las tarjetas de crédito para transacciones y pagos a plazos en este cluster. Esto se debe a que el cliente tiene un saldo relativamente pequeño, la frecuencia del saldo cambia rara vez y los pagos a plazos son muy bajos. Además, un límite de crédito bajo también muestra que los clientes rara vez o casi nunca usan las tarjetas de crédito para procesar transacciones de crédito, y los clientes en este cluster también rara vez hacen adelantos en efectivo. Por lo tanto, se puede suponer que los clientes utilizan las tarjetas de crédito para procesos de adelanto de efectivo solo con la suficiente frecuencia. Además, el bajo saldo permite que los clientes de este cluster sean estudiantes o nuevos usuarios que utilizan las tarjetas de crédito en este banco.
- **Cluster 3 (Usuarios de Pagos a Plazos):** En este cluster, los clientes utilizan las tarjetas de crédito específicamente para pagos a plazos. Esto se debe al nivel relativamente alto de transacciones que utilizan pagos a plazos en este cluster. Además, los clientes en este cluster a menudo realizan transacciones con montos muy grandes por transacción y la frecuencia y transacciones de adelantos en efectivo son muy pequeñas. Los clientes en este cluster rara vez hacen pagos y adelantos en efectivo y tienen una frecuencia de adelantos en efectivo y un monto de pagos relativamente pequeños. Se puede concluir que los clientes en este cluster son muy adecuados para tarjetas de crédito específicamente para necesidades de pagos a plazos.
- **Cluster 4 (Usuarios de Adelantos en Efectivo/Retiros):** Los clientes en este cluster tienen saldos altos, la frecuencia de los saldos siempre cambia y la frecuencia de adelantos en efectivo y pagos anticipados es alta. Además, los clientes en este cluster tienen las tasas de interés más bajas en comparación con los otros clusters y tienen el segundo límite de crédito más alto y los pagos de los cuatro clusters. Sin embargo, los usuarios de tarjetas de crédito en este cluster rara vez hacen pagos a plazos o compras únicas y tienen la tercera antigüedad más alta de los cuatro clusters. Por lo tanto, se puede concluir que los clientes en este cluster solo usan tarjetas de crédito para la necesidad de retirar dinero o hacer adelantos en efectivo.

Las siguientes son visualizaciones sobre 2 variables para cada cluster:

In [302]:

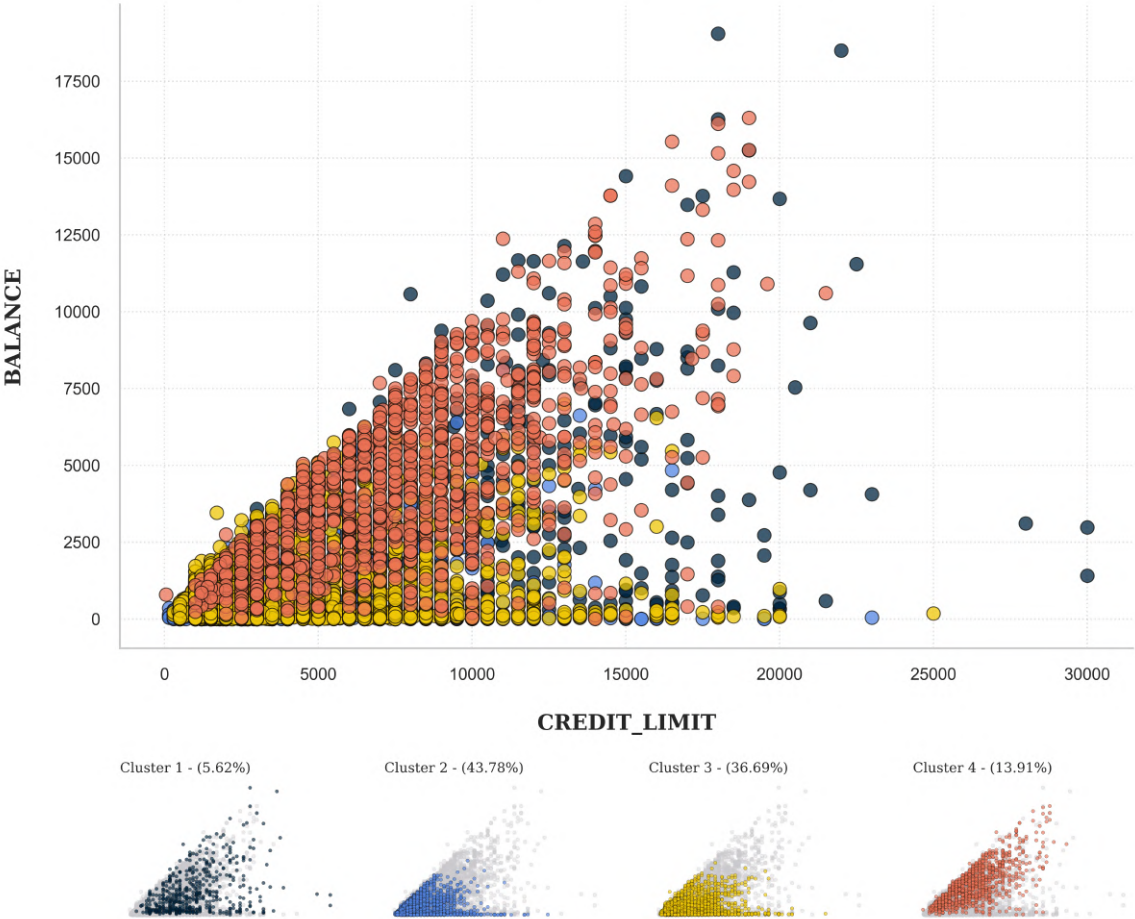
```

1 # --- Cluster Visualization 1: Variables ---
2 scatter_style=dict(linewidth=0.65, edgecolor='#100C07', alpha=0.75)
3 sub_scatter_style_color=dict(s=5, alpha=0.65, linewidth=0.15, zorder=10, edgecolor='#
4 sub_scatter_style_grey=dict(s=5, alpha=0.3, linewidth=0.7, zorder=5, color='#CAC9CD')
5 grid_style=dict(alpha=0.3, color='#9B9A9C', linestyle='dotted', zorder=1)
6 xy_label=dict(fontweight='bold', fontsize=14, fontfamily='serif')
7 suptitle=dict(fontsize=20, fontweight='heavy', fontfamily='serif')
8 title=dict(fontsize=14, fontfamily='serif')
9 color_pallete=['#002642', '#5386E4', '#F0C808', '#EC7357']
10 sub_axes=[None] * 4
11
12 # --- Cluster Visualization 1: Data Frame ---
13 df_cv1 = df[['CREDIT_LIMIT', 'BALANCE', 'cluster_result']]
14 cluster_result = sorted(df_cv1['cluster_result'].unique())
15
16 # --- Cluster Visualization 1: Settings ---
17 fig = plt.figure(figsize=(12, 16))
18 gs = fig.add_gridspec(4, 4)
19 ax = fig.add_subplot(gs[:4, :])
20 ax.set_aspect(1)
21
22 # --- Cluster Visualization 1: Main Scatter Plot ---
23 for x in range(len(cluster_result)):
24     df_cv1_x = df_cv1[df_cv1['cluster_result']==cluster_result[x]]
25
26     ax.scatter(df_cv1_x['CREDIT_LIMIT'], df_cv1_x['BALANCE'], s=80, color=color_palle
27     ax.set_title('Clusters 1 and 4 have the highest balance and credit limit compared
28     ax.set_xlabel('\nCREDIT_LIMIT', **xy_label)
29     ax.set_ylabel('BALANCE\n', **xy_label)
30     ax.grid(axis='y', which='major', **grid_style)
31     ax.grid(axis='x', which='major', **grid_style)
32     for spine in ax.spines.values():
33         spine.set_color('None')
34     for spine in ['bottom', 'left']:
35         ax.spines[spine].set_visible(True)
36         ax.spines[spine].set_color('#CAC9CD')
37     plt.xticks(fontsize=11)
38     plt.yticks(fontsize=11)
39
40 # --- Cluster Visualization 1: Sub Plots ---
41 for idx, clstr in enumerate(cluster_result):
42     sub_axes[idx] = fig.add_subplot(gs[3, idx], aspect=1)
43
44     sub_axes[idx].scatter(df_cv1[df_cv1['cluster_result']!=clstr]['CREDIT_LIMIT'], df
45     sub_axes[idx].scatter(df_cv1[df_cv1['cluster_result']==clstr]['CREDIT_LIMIT'], df
46
47     cnt = round((df_cv1['cluster_result']==clstr).sum()/8950*100, 2)
48     sub_axes[idx].set_title(f'{clstr} - ({cnt}%)', loc='left', fontsize=9, fontfamily
49     sub_axes[idx].set_xticks([])
50     sub_axes[idx].set_yticks([])
51     for spine in sub_axes[idx].spines.values():
52         spine.set_color('None')
53
54 # --- Cluster Visualization 1: Title ---
55 plt.suptitle('Credit Limit vs. Balance based on Clusters', x=0.123, y=0.73, ha='left'
56

```

Credit Limit vs. Balance based on Clusters

Clusters 1 and 4 have the highest balance and credit limit compared to other clusters.



In [303]:

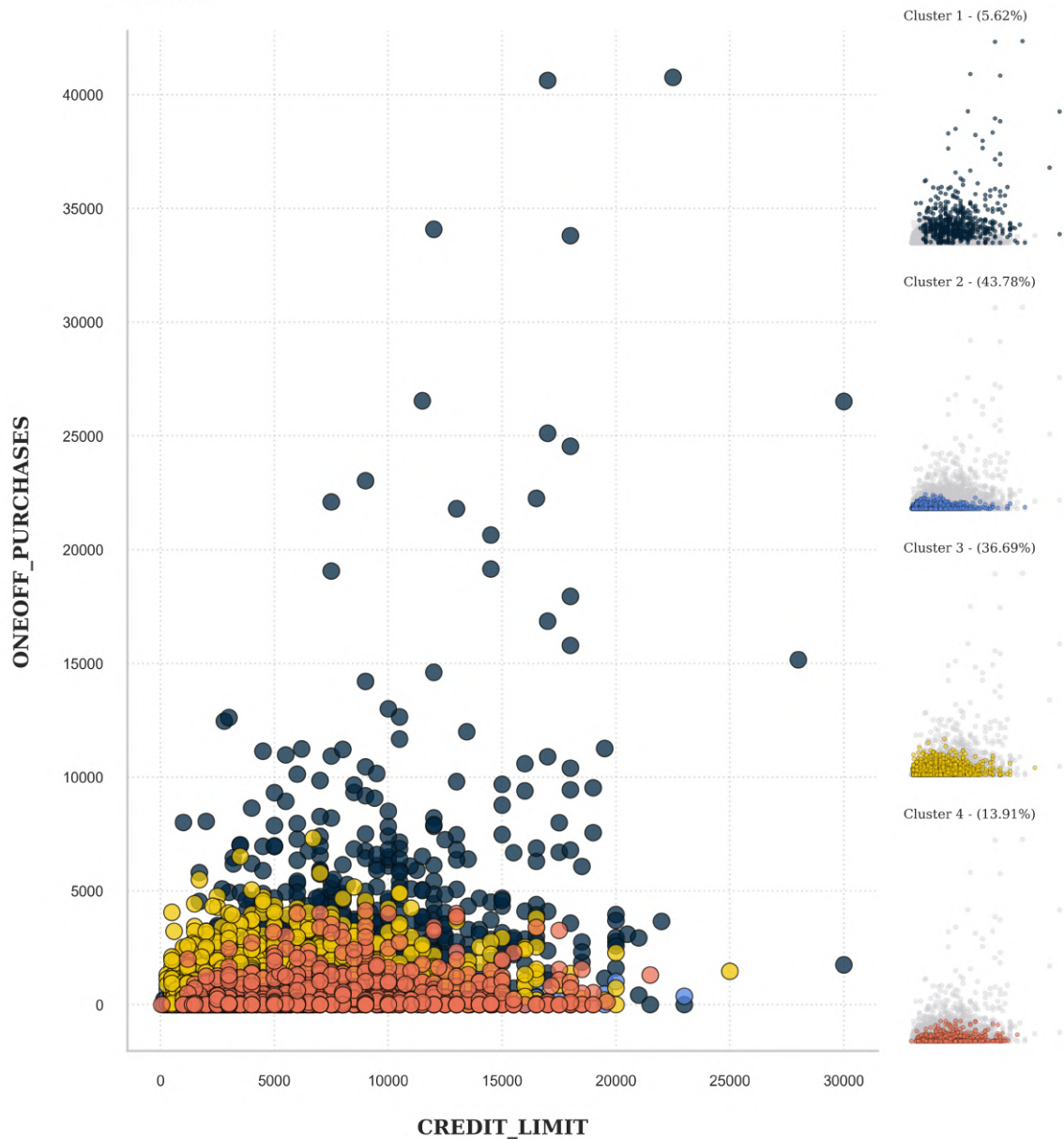
```

1 # --- Cluster Visualization 2: Variables ---
2 scatter_style=dict(linewidth=0.65, edgecolor='#100C07', alpha=0.75)
3 sub_scatter_style_color=dict(s=5, alpha=0.65, linewidth=0.15, zorder=10, edgecolor='#
4 sub_scatter_style_grey=dict(s=5, alpha=0.3, linewidth=0.7, zorder=5, color='#CAC9CD')
5 grid_style=dict(alpha=0.3, color='#9B9A9C', linestyle='dotted', zorder=1)
6 xy_label=dict(fontsize=11, fontweight='bold', fontfamily='serif')
7 suptitle=dict(fontsize=14, fontweight='heavy', fontfamily='serif')
8 title=dict(fontsize=11, fontfamily='serif')
9 color_pallete=['#002642', '#5386E4', '#F0C808', '#EC7357']
10 sub_axes=[None] * 4
11
12 # --- Cluster Visualization 2: Data Frame ---
13 df_cv2 = df[['CREDIT_LIMIT', 'ONEOFF_PURCHASES', 'cluster_result']]
14 cluster_result = sorted(df_cv1['cluster_result'].unique())
15
16 # --- Cluster Visualization 2: Settings ---
17 fig = plt.figure(figsize=(12, 10))
18 gs = fig.add_gridspec(4, 4)
19 ax = fig.add_subplot(gs[:4, :4])
20 ax.set_aspect(1)
21
22 # --- Cluster Visualization 2: Main Scatter Plot ---
23 for x in range(len(cluster_result)):
24     df_cv2_x = df_cv2[df_cv2['cluster_result']==cluster_result[x]]
25
26     ax.scatter(df_cv2_x['CREDIT_LIMIT'], df_cv2_x['ONEOFF_PURCHASES'], s=80, color=cc
27     ax.set_title('There is no correlation between the one-off purchase amount and the
28     ax.set_xlabel('\nCREDIT_LIMIT', **xy_label)
29     ax.set_ylabel('ONEOFF_PURCHASES\n', **xy_label)
30     ax.grid(axis='y', which='major', **grid_style)
31     ax.grid(axis='x', which='major', **grid_style)
32     for spine in ax.spines.values():
33         spine.set_color('None')
34     for spine in ['bottom', 'left']:
35         ax.spines[spine].set_visible(True)
36         ax.spines[spine].set_color('#CAC9CD')
37     plt.xticks(fontsize=8)
38     plt.yticks(fontsize=8)
39
40 # --- Cluster Visualization 2: Sub Plots ---
41 for idx, clstr in enumerate(cluster_result):
42     sub_axes[idx] = fig.add_subplot(gs[idx, 3], aspect=1)
43
44     sub_axes[idx].scatter(df_cv2[df_cv2['cluster_result']!=clstr]['CREDIT_LIMIT'], df
45     sub_axes[idx].scatter(df_cv2[df_cv2['cluster_result']==clstr]['CREDIT_LIMIT'], df
46
47     cnt = round((df_cv2['cluster_result']==clstr).sum()/8950*100, 2)
48     sub_axes[idx].set_title(f'{clstr} - ({cnt}%)', loc='left', fontsize=7, fontfamily
49     sub_axes[idx].set_xticks([])
50     sub_axes[idx].set_yticks([])
51     for spine in sub_axes[idx].spines.values():
52         spine.set_color('None')
53
54 # --- Cluster Visualization 2: Title ---
55 plt.suptitle('One-off Purchase vs. Credit Limit based on Clusters', x=0.275, y=0.96,
56

```

One-off Purchase vs. Credit Limit based on Clusters

There is no correlation between the one-off purchase amount and the credit limit obtained.



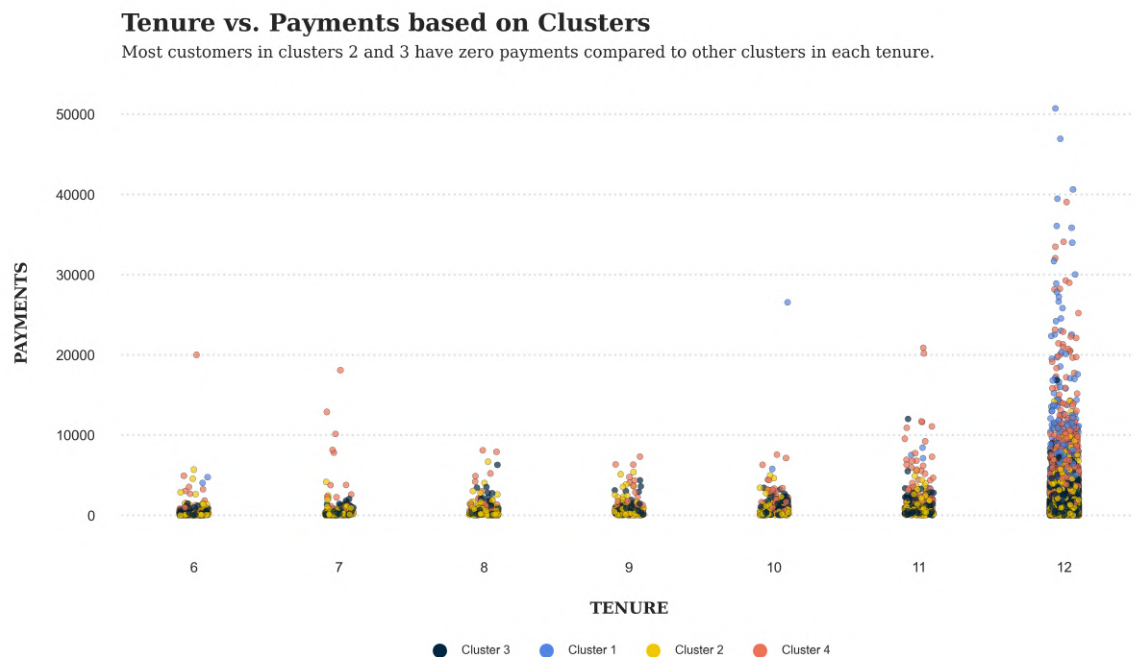
Las compras al contado no tiene relacion (aparente) con el límite de crédito adicional obtenido por el usuario. En la figura de arriba, como se mencionó anteriormente, se puede ver que el grupo 1 tiene un cliente con la mayor cantidad de compra para una transacción.

In [304]:

```

1 # --- Cluster Visualization 3: Data Frame ---
2 df_cv3 = df[['TENURE', 'PAYMENTS', 'cluster_result']]
3
4 # --- Cluster Visualization 3: Variables ---
5 color_pallete = ['#002642', '#5386E4', '#F0C808', '#EC7357']
6 supitle = dict(fontsize=12, ha='left', fontweight='heavy', fontfamily='serif')
7 title = dict(fontsize=8, loc='left', fontfamily='serif')
8 cluster_result = sorted(df_cv3['cluster_result'].unique())
9 stripplot_style = dict(edgecolor='#100C07', s=3, linewidth=0.15, alpha=0.7, palette=c
10 legend_style = dict(ncol=5, borderpad=3, frameon=False, fontsize=6, title=None)
11 xy_label = dict(fontweight='bold', fontsize=8, fontfamily='serif')
12 grid_style = dict(alpha=0.3, color='#9B9A9C', linestyle='dotted', zorder=1)
13
14 # --- Cluster Visualization 3: Visuals ---
15 stplot=sns.stripplot(data=df_cv3, x='TENURE', y='PAYMENTS', hue='cluster_result', **s
16 sns.move_legend(stplot, 'upper center', bbox_to_anchor=(0.5, -0.15), **legend_style)
17 sns.despine(top=True, right=True, left=True, bottom=True)
18 plt.suptitle('Tenure vs. Payments based on Clusters', x=0.125, y=1.01, **suptitle)
19 plt.title('Most customers in clusters 2 and 3 have zero payments compared to other cl
20 plt.xlabel('\nTENURE', **xy_label)
21 plt.ylabel('PAYMENTS\n', **xy_label)
22 plt.xticks(fontsize=7)
23 plt.yticks(fontsize=7)
24 plt.grid(axis='x', alpha=0)
25 plt.grid(axis='y', **grid_style)
26
27 plt.gcf().set_size_inches(9, 4)
28 plt.show();

```



La mayoría de los clientes en los grupos 2 y 3 no tienen pagos en comparación con otros grupos en cada periodo de tiempo. Como se mencionó anteriormente, se puede observar que la mayoría de los clientes tienden a elegir un plazo de 12 meses.

In [305]:

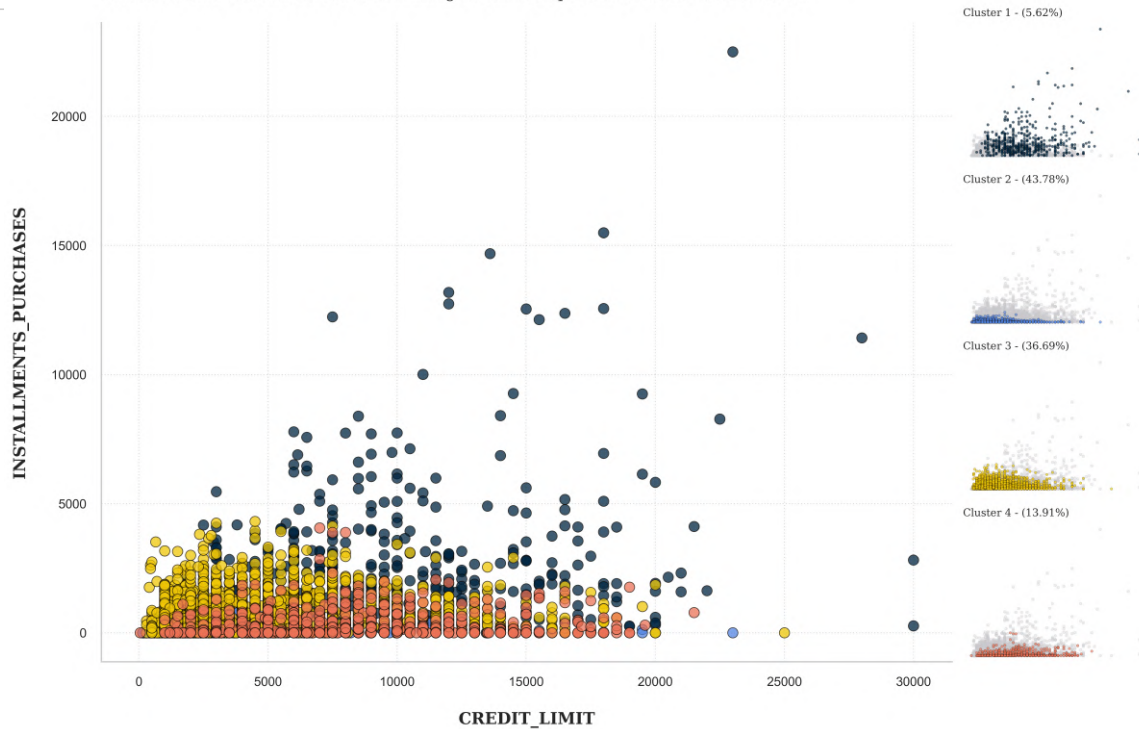
```

1 # --- Cluster Visualization 4: Data Frame ---
2 df_cv4 = df[['INSTALLMENTS_PURCHASES', 'CREDIT_LIMIT', 'cluster_result']]
3
4 # --- Cluster Visualization 4: Variables ---
5 cluster_result = sorted(df_cv4['cluster_result'].unique())
6 scatter_style=dict(linewidth=0.65, edgecolor='#100C07', alpha=0.75)
7 sub_scatter_style_color=dict(s=5, alpha=0.65, linewidth=0.15, zorder=10, edgecolor='#
8 sub_scatter_style_grey=dict(s=5, alpha=0.3, linewidth=0.7, zorder=5, color='#CAC9CD')
9 grid_style=dict(alpha=0.3, color='#9B9A9C', linestyle='dotted', zorder=1)
10 xy_label=dict(fontsize=14, fontweight='bold', fontfamily='serif')
11 supitle=dict(fontsize=20, fontweight='heavy', fontfamily='serif')
12 title=dict(fontsize=14, fontfamily='serif')
13 xy_label=dict(fontweight='bold', fontsize=14, fontfamily='serif')
14 grid_style=dict(alpha=0.3, color='#9B9A9C', linestyle='dotted', zorder=1)
15 color_pallete=['#002642', '#5386E4', '#F0C808', '#EC7357']
16 sub_axes=[None] * 4
17
18 # --- Cluster Visualization 4: Settings ---
19 fig = plt.figure(figsize=(21, 10))
20 gs = fig.add_gridspec(4, 4)
21 ax = fig.add_subplot(gs[:4, :4])
22 ax.set_aspect(1)
23
24 # --- Cluster Visualization 4: Main Scatter Plot ---
25 for x in range(len(cluster_result)):
26     df_cv4_x = df_cv4[df_cv4['cluster_result']==cluster_result[x]]
27
28     ax.scatter(df_cv4_x['CREDIT_LIMIT'], df_cv4_x['INSTALLMENTS_PURCHASES'], s=80, co
29 ax.set_title('Clusters 1 and 3 are more active in making installment purchases th
30 ax.set_xlabel('\nCREDIT_LIMIT', **xy_label)
31 ax.set_ylabel('INSTALLMENTS_PURCHASES\n', **xy_label)
32 ax.grid(axis='y', which='major', **grid_style)
33 ax.grid(axis='x', which='major', **grid_style)
34 for spine in ax.spines.values():
35     spine.set_color('None')
36 for spine in ['bottom', 'left']:
37     ax.spines[spine].set_visible(True)
38     ax.spines[spine].set_color('#CAC9CD')
39 plt.xticks(fontsize=11)
40 plt.yticks(fontsize=11)
41
42 # --- Cluster Visualization 4: Sub Plots ---
43 for idx, clstr in enumerate(cluster_result):
44     sub_axes[idx] = fig.add_subplot(gs[idx, 3], aspect=1)
45
46     sub_axes[idx].scatter(df_cv4[df_cv4['cluster_result']!=clstr]['CREDIT_LIMIT'], df
47     sub_axes[idx].scatter(df_cv4[df_cv4['cluster_result']==clstr]['CREDIT_LIMIT'], df
48
49     cnt = round((df_cv4['cluster_result']==clstr).sum()/8950*100, 2)
50     sub_axes[idx].set_title(f'{clstr} - ({cnt}%)', loc='left', fontsize=9, fontfamily
51     sub_axes[idx].set_xticks([])
52     sub_axes[idx].set_yticks([])
53     for spine in sub_axes[idx].spines.values():
54         spine.set_color('None')
55
56 # --- Cluster Visualization 4: Title ---
57 plt.suptitle('Installments Purchases vs. Credit Limit based on Clusters', x=0.268, y=
58

```

Installments Purchases vs. Credit Limit based on Clusters

Clusters 1 and 3 are more active in making installment purchases than other clusters.



It can be seen that clusters 1 and 3 have more installment purchases than clusters 2 and 4. However, it can also be seen that a large number of installment purchases are not correlated with the credit limit increase.

8.2 | Strategy Suggestions 💡

Basándonos en los resultados del perfil anterior y en el análisis realizado para cada uno de los clusters, aquí hay algunas sugerencias para estrategias:

- **Los clientes del cluster 1** pueden convertirse en el objetivo principal para el marketing de tarjetas de crédito. **Esto se debe a que los clientes en este cluster son muy activos en el uso de tarjetas de crédito y tienen la antigüedad y los límites de crédito más altos en comparación con los otros clusters.** Al enfocar el marketing en este cluster, los bancos pueden aumentar sus ganancias al usar más tarjetas de crédito y reducir los costos de marketing incurridos. Los bancos pueden ofrecer beneficios o recompensas por el uso de tarjetas de crédito para atraer a los clientes a usar las tarjetas de crédito con más frecuencia.
- Para las tarjetas de crédito para pagos a plazos, los bancos pueden centrar su marketing en los clientes del cluster 3. Esto se debe a que **los clientes en el cluster 3 son más propensos a realizar transacciones con tarjeta de crédito para fines de pagos a plazos.** Los bancos pueden ofrecer programas de pagos a plazos con intereses bajos o del 0% que se puedan utilizar para diversas necesidades de pagos a plazos a los clientes en este cluster para atraer a los clientes a utilizar las tarjetas de crédito. **Los requisitos de pagos a plazos que se pueden ofrecer pueden ser en forma de pagos de viajes, electrodomésticos, dispositivos electrónicos, smartphones o ciertas marcas que son más demandadas por el público.**
- **Para el cluster 2**, los bancos pueden ofrecer tarjetas de crédito especiales para principiantes o estudiantes (tarjetas de nivel de entrada) que puedan carecer de un perfil crediticio extenso. **Esta tarjeta de crédito se puede utilizar para construir su crédito y aprender a utilizar la tarjeta de crédito de manera responsable.** Esta tarjeta puede incluir varias características, como **exención de tarifas, recompensas** por establecer una rutina de pagos puntuales, **barreras bajas** para convertirse en titular de la tarjeta y tasas de interés indulgentes. Además, los bancos pueden **ofrecer**

oportunidades para actualizar a nuevos productos y mejores términos y condiciones si el cliente paga consistentemente las facturas a tiempo. Además, los bancos pueden ofrecer recompensas de registro para que los clientes que no son usuarios de la tarjeta de crédito del banco puedan interesarse en registrarse.

- Dado que los clientes en **el cluster 4 tienden a realizar adelantos en efectivo**, los bancos pueden ofrecer tarjetas de crédito especiales con diversos beneficios. Estos beneficios pueden ser en forma de bajos o ningún cargo por adelantos en efectivo o administrativos, bajos intereses, antigüedad relativamente alta, etc. Además, los bancos también **pueden ofrecer programas bancarios distintos de las tarjetas de crédito, como programas de pay-later con colaboraciones de terceros o préstamos personales proporcionados por los bancos.**

In []:

1	
---	--