

## ✓ Librerías

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import requests
6 import os
7 import zipfile
8 from datetime import datetime
9 !pip install tabula-py
10 import tabula
11 import re
```

```
Requirement already satisfied: tabula-py in /usr/local/lib/python3.12/dist-packa
Requirement already satisfied: pandas>=0.25.3 in /usr/local/lib/python3.12/dist-
Requirement already satisfied: numpy>1.24.4 in /usr/local/lib/python3.12/dist-pa
Requirement already satisfied: distro in /usr/local/lib/python3.12/dist-packages
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.1
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.12/dist-pa
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.12/dist-
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.12/dist-packag
```

## ✓ Recopilación de Datos

```
1 URL_DESCARGA = "https://estadisticas.seps.gob.ec/?sdm_process_download=1&dc
2 CARPETA_BASE = "seps_descargas"
3 CARPETA_SEG1 = "seps_segmento1"
4
5 os.makedirs(CARPETA_BASE, exist_ok=True)
6 os.makedirs(CARPETA_SEG1, exist_ok=True)
7
8 def descargar_archivo():
9     fecha = datetime.now().strftime("%Y%m%d")
10    ruta_archivo = f"{CARPETA_BASE}/seps_{fecha}.zip"
11
12    print("Descargando archivo...")
13    r = requests.get(URL_DESCARGA)
14    with open(ruta_archivo, "wb") as f:
15        f.write(r.content)
16
17    print(f"Archivo descargado: {ruta_archivo}")
18    return ruta_archivo
19
```

```

20 def extraer_segmento1(ruta_archivo):
21     fecha = datetime.now().strftime("%Y%m%d")
22     carpeta_temp = f"{CARPETA_BASE}/extraido_{fecha}"
23     os.makedirs(carpeta_temp, exist_ok=True)
24
25     # Extraer ZIP
26     with zipfile.ZipFile(ruta_archivo, "r") as z:
27         z.extractall(carpeta_temp)
28
29     print("Archivos extraídos. Buscando archivo Segmento 1...")
30
31     # Buscar el archivo Segmento 1
32     archivo_segmento1 = None
33     for root, dirs, files in os.walk(carpeta_temp):
34         for file in files:
35             if "segmento" in file.lower() and "1" in file.lower() and file.
36                 archivo_segmento1 = os.path.join(root, file)
37
38     if archivo_segmento1 is None:
39         raise FileNotFoundError("No se encontró el archivo Segmento 1 dentro de la carpeta")
40
41     # Copiar a carpeta final
42     destino = os.path.join(CARPETA_SEG1, f"{fecha}_segmento1.xlsm")
43     with open(archivo_segmento1, "rb") as f_origen:
44         with open(destino, "wb") as f_destino:
45             f_destino.write(f_origen.read())
46
47     print(f"Archivo Segmento 1 listo: {destino}")
48     return destino
49
50
51
52 zipfile_path = descargar_archivo()
53 ruta_segmento1 = extraer_segmento1(zipfile_path)
54

```

Descargando archivo...

Archivo descargado: seps\_descargas/seps\_20251204.zip

Archivos extraídos. Buscando archivo Segmento 1...

Archivo Segmento 1 listo: seps\_segmento1/20251204\_segmento1.xlsm

```

1 URL = "https://www.seps.gob.ec/?sdm_process_download=1&download_id=28370"
2 NOMBRE_ARCHIVO = "Calificacion_de_Riesgo.pdf"
3
4 print(f"Descargando PDF desde:\n{URL}\n")
5
6 resp = requests.get(URL)
7 resp.raise_for_status() # Si hay error, detiene el programa
8
9 with open(NOMBRE_ARCHIVO, "wb") as f:
10     f.write(resp.content)

```

```
11
12 print(f"PDF guardado como: {NOMBRE_ARCHIVO}")
13
14
15 pdf_path = "/content/Calificacion_de_Riesgo.pdf"
16
17 tablas = tabula.read_pdf(pdf_path, pages="all", lattice=True, multiple_tabl
18 df_ranking = pd.concat(tablas, ignore_index=True)
19
20
```

Descargando PDF desde:

[https://www.seps.gob.ec/?sdm\\_process\\_download=1&download\\_id=28370](https://www.seps.gob.ec/?sdm_process_download=1&download_id=28370)

PDF guardado como: Calificacion\_de\_Riesgo.pdf

## ✓ Limpieza y modificación de datos

## ✓ ASPECTOS TECNICOS

```
1 df_aspectos_tecnicos = pd.read_excel(ruta_segmento1, sheet_name='3. RESULTA
2 print("3. 'RESULTADOS DEL EJERCICIO' cargados. Primeras 5 filas:")
3 df_aspectos_tecnicos.head()
```

3. 'RESULTADOS DEL EJERCICIO' cargados. Primeras 5 filas:  
/usr/local/lib/python3.12/dist-packages/openpyxl/worksheet/\_reader.py:329: UserWarning(msg)

Unnamed: 0			15 DE ABRIL LTDA	23 DE JULIO LIMITADA	29 DE OCTUBRE LTDA	9 DE OCTUBRE LTDA	JAN DE LE
0	5	Ingresos	13311454.95	48247840.25	1.201518e+08	13563505.43	163
1	51	Intereses y descuentos ganados	10756423.83	41868717.33	1.054787e+08	12615572.14	152
2	(-) 41	Intereses causados	5673237.15	25450929.49	6.340963e+07	6952973.27	92
3	NaN	MARGEN NETO DE INTERESES	5083186.68	16417787.84	4.206908e+07	5662598.87	59
4	(+) 52	Comisiones ganadas	0.00	159313.57	5.582783e+06	0.00	3

5 rows × 46 columns

ESTADOS FINANCIEROS

```
1 df_estado_financiero = pd.read_excel(ruta_segmento1, sheet_name='2. ESTADO
2 print("Hoja '2. ESTADO FINANCIERO'. Primeras 5 filas:")
3 df_estado_financiero.head()
```

Hoja '2. ESTADO FINANCIERO'. Primeras 5 filas:  
/usr/local/lib/python3.12/dist-packages/openpyxl/worksheet/\_reader.py:329: UserWarning: warn(msg)

	COD CONTABLE	Nombre de Cuenta	TIPO*	GRUPO**	15 DE ABRIL LTDA	23 DE JULIO LIMITADA	29 OCTUB LT
0	1	ACTIVO	1	1	1.293271e+08	4.437386e+08	1.144214e+
1	11	FONDOS DISPONIBLES	1	2	1.442083e+07	7.201578e+07	1.436635e+
2	1101	Caja	1	4	1.193649e+06	6.805766e+06	8.360886e+
3	110105	Efectivo	1	6	1.193449e+06	6.794486e+06	8.316997e+
4	110110	Caja chica	1	6	2.000000e+02	1.128000e+04	4.388895e+

5 rows × 48 columns

> **INDICADORES FINANCIEROS**

↳ 7 cells hidden

✓ **Ranking**

Quitar todos los na, saltos de linea y espacios

```
1 df_ranking = df_ranking.dropna(how="all")
2 df_ranking = df_ranking.applymap(lambda x: str(x).replace("\n", " ").strip())
3
4 df_ranking
```

/tmp/ipython-input-15491972.py:2: FutureWarning: DataFrame.applymap has been deprecated  
df\_ranking = df\_ranking.applymap(lambda x: str(x).replace("\n", " ").strip())

Unnamed: 0	Unnamed: 1	Unnamed: 2	Unnamed: 3	Unnamed: 4
0	No.	RUC		INSTITUCIÓN FINANCIERA
1	1	1768168480001	CORPORACIÓN NACIONAL DE FINANZAS	POPULARES Y RURALES
3	No.	RUC		INSTITUCIÓN FINANCIERA
4	2	1890012015001	ASOCIACIÓN MUTUALISTA DE AHO	RO Y CRÉDITO PARA EL DESARROLLO
5	3	0190006247001	ASOCIACIÓN MUTUALISTA DE AHO	RO Y CRÉDITO PARA EL DESARROLLO
...	...	...		
61	58	0790024656001	COOPERATIVA DE AHO	RO Y CRÉDITO SANTA ROSA LIMITADA
62	59	0490002669001	COOPERATIVA DE AHO	RO Y CRÉDITO TULCÁN LIMITADA
63	60	1190068389001	COOPERATIVA DE AHO	RO Y CRÉDITO VICENTINAMANUELA
64	61	0591711164001	COOPERATIVA DE AHO	RO Y CRÉDITO VIRGEN DEL CARMEN
65	62	0391005664001	COOPERATIVA DE AHO	RO Y CRÉDITO YUYAY

65 rows × 10 columns

Next steps:

[Generate code with df\\_ranking](#)

[New interactive sheet](#)

Modificar para los encabezados

```
1 header = df_ranking.iloc[0].tolist()      # fila 0 = encabezado
2 df = df_ranking[1:].copy()                # borrar fila de encabezado
3 df.columns = header                      # asignar encabezado
4
5 df
```

	No.	RUC	INSTITUCIÓN FINANCIERA
1	1	1768168480001	CORPORACIÓN NACIONAL DE FINANZAS \r POPULARES Y SOL...
3	No.	RUC	INSTITUCIÓN FINANCIERA
4	2	1890012015001	ASOCIACION MUTUALISTA DE AHORRO Y \r CREDITO PARA LA...
5	3	0190006247001	ASOCIACION MUTUALISTA DE AHORRO Y \r CREDITO PARA LA...
6	4	1090056286001	ASOCIACION MUTUALISTA DE AHORRO Y \r CREDITO PARA LA...
...	...	...	...
61	58	0790024656001	COOPERATIVA DE AHORRO Y CRÉDITOS SANTA \r ROSA LIMITADA.
62	59	0490002669001	COOPERATIVA DE AHORRO Y CRÉDITO TULCÁN \r LIMITADA.
63	60	1190068389001	COOPERATIVA DE AHORRO Y CRÉDITO \r VICENTIN AMANUELES...
64	61	0591711164001	COOPERATIVA DE AHORRO Y CRÉDITO VIRGEN \r DEL CISNE

Next steps:

[Generate code with df](#)[New interactive sheet](#)

Limpiar los nombres de las columnas para que sea mas facil manejarlos luego

```
1 df.columns = [c.replace("\r", " ").replace(" ", " ").strip() for c in df.c
2
3 rename_map = {
4     "N": "N",
5     "RUC": "RUC",
6     "INSTITUCIÓN FINANCIERA": "Institucion",
7     "SEGMENTO": "Segmento",
8     "FIRMA CALIFICADORA DE RIESGO": "Firma",
9     "AL 30 DE JUNIO 2024": "Jun2024",
10    "AL 30 DE SEPTIEMBRE 2024": "Sep2024",
11    "AL 31 DE DICIEMBRE 2024": "Dic2024",
12    "AL 31 DE MARZO 2025": "Mar2025",
13    "AL 30 DE JUNIO 2025": "Jun2025"
14 }
15
16 df = df.rename(columns=rename_map)
17 df
```



No.		RUC		Institucion
1	1	1768168480001	CORPORACIÓN NACIONAL DE FINANZAS	POPULARES Y SOLIDARIAS
3	No.	RUC	INSTITUCIÓN FINANCIERA	
4	2	1890012015001	ASOCIACION MUTUALISTA DE AHOORRO Y	CREDITO PARA LA POBLACION
5	3	0190006247001	ASOCIACION MUTUALISTA DE AHOORRO Y	CREDITO PARA LA POBLACION
6	4	1090056286001	ASOCIACION MUTUALISTA DE AHOORRO Y	CREDITO PARA LA POBLACION
...	...	...	...	
61	58	0790024656001	COOPERATIVA DE AHOORRO Y CREDITO	SANTA ROSA LIMITADA
62	59	0490002669001	COOPERATIVA DE AHOORRO Y CREDITO	TULCAN LIMITADA
63	60	1190068389001	COOPERATIVA DE AHOORRO Y CREDITO	VICENTIN AMANUELES
64	61	0591711164001	COOPERATIVA DE AHOORRO Y CREDITO	VIRGEN DEL CISNE

Next steps: [Generate code with df](#) [New interactive sheet](#)

1 df

No .		RUC		Institucion
1	1	1768168480001	CORPORACIÓN NACIONAL DE FINANZAS	POPULARES Y SOLIDARIAS
3	No.	RUC	INSTITUCIÓN FINANCIERA	
4	2	1890012015001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA LA	COMUNIDAD
5	3	0190006247001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA LA	COMUNIDAD
6	4	1090056286001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA LA	COMUNIDAD
...	...	...	...	
61	58	0790024656001	COOPERATIVA DE AHORRO Y CREDITO SANTA ROSA	LIMITADA
62	59	0490002669001	COOPERATIVA DE AHORRO Y CREDITO TULCAN	LIMITADA
63	60	1190068389001	COOPERATIVA DE AHORRO Y CREDITO VICENTINAMANUELES	...
64	61	0591711164001	COOPERATIVA DE AHORRO Y CREDITO VIRGEN DEL CISNE	...

Next steps: [Generate code with df](#) [New interactive sheet](#)

Limpiar columnas

```
1 def limpiar(s):
2     if pd.isna(s):
3         return None
4     s = str(s)
```

```
5     s = s.replace("-", "-").replace("-", "-") # Normalizar guiones
6     s = s.upper()
7     s = s.replace("`", "'")
8     s = re.sub(r"\s*\|\s*", " | ", s)          # Normalizar separador |
9     s = re.sub(r"\s*-\s*", "-", s)            # Normalizar guiones con espa
10    s = re.sub(r"\s{2,}", " ", s)              # Quitar múltiples espacios
11    return s.strip()
12
13 df = df.applymap(limpiar)
14 df
```

```
/tmp/ipython-input-2218334321.py:13: FutureWarning: DataFrame.applymap has been
df = df.applymap(limpiar)
```

No.		RUC	Institucion
1	1	1768168480001	CORPORACIÓN NACIONAL DE FINANZAS \r POPULARES Y SOL...
3	NO.	RUC	INSTITUCIÓN FINANCIERA
4	2	1890012015001	ASOCIACION MUTUALISTA DE AHOORRO Y \r CREDITO PARA LA...
5	3	0190006247001	ASOCIACION MUTUALISTA DE AHOORRO Y \r CREDITO PARA LA...
6	4	1090056286001	ASOCIACION MUTUALISTA DE AHOORRO Y \r CREDITO PARA LA...
...	...	...	...
61	58	0790024656001	COOPERATIVA DE AHOORRO Y CRÉDITO SANTA \r ROSA LIMITADA.
62	59	0490002669001	COOPERATIVA DE AHOORRO Y CRÉDITO TULCÁN \r LIMITADA.
63	60	1190068389001	COOPERATIVA DE AHOORRO Y CRÉDITO \r VICENTIN AMANUELES...
64	61	0591711164001	COOPERATIVA DE AHOORRO Y CRÉDITO VIRGEN \r DEL CISNE

Next steps:

Generate code with df

New interactive sheet

Filtrar solo para segmento 1

```
1 #Normalizar la columna 'Segmento' para extraer solo el número
2 if 'Segmento' not in df.columns:
```

```

3     raise ValueError("No se encontró la columna 'Segmento' en df.columns: ")
4
5 #Extrae el primer número que aparezca en la celda (ej: 'SEGMENTO 1', 'SEGME
6 df['Segmento'] = df['Segmento'].astype(str).str.extract(r'(\d+)')[0]
7 # Si prefieres mantener como int:
8 df['Segmento'] = df['Segmento'].astype(float).astype('Int64') # permite NA
9
10 # Limpiar columnas de calificación (quitar asteriscos y quedarnos con el tc
11 rating_cols = [c for c in df.columns if re.search(r'Jun|Sep|Dic|Mar|Jun', c
12 for c in rating_cols:
13     # nos quedamos con la última "palabra" (p.ej. 'AAA- *****' -> 'AAA-')
14     df[c] = (df[c].astype(str)
15             .str.split()
16             .str[0]                # o str[-1] según el patrón; usar str[0]
17             .str.replace(r'\*+', '', regex=True)
18             .str.strip()
19             .replace('', pd.NA))
20
21 # filtrar Segmento == 1
22 df_ranking = df[df['Segmento'] == 1].copy()
23
24 print("Shape total limpio:", df.shape)
25 print("Shape df_ranking (Segmento 1):", df_ranking.shape)
26 display(df_ranking.head(10))

```

Shape total limpio: (64, 10)  
 Shape df\_ranking (Segmento 1): (49, 10)

	No.	RUC	Institucion	Segmen
1	1	1768168480001	CORPORACIÓN NACIONAL DE FINANZAS POPULARES Y SOL...	
4	2	1890012015001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA LA...	
5	3	0190006247001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA LA...	
6	4	1090056286001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA LA...	

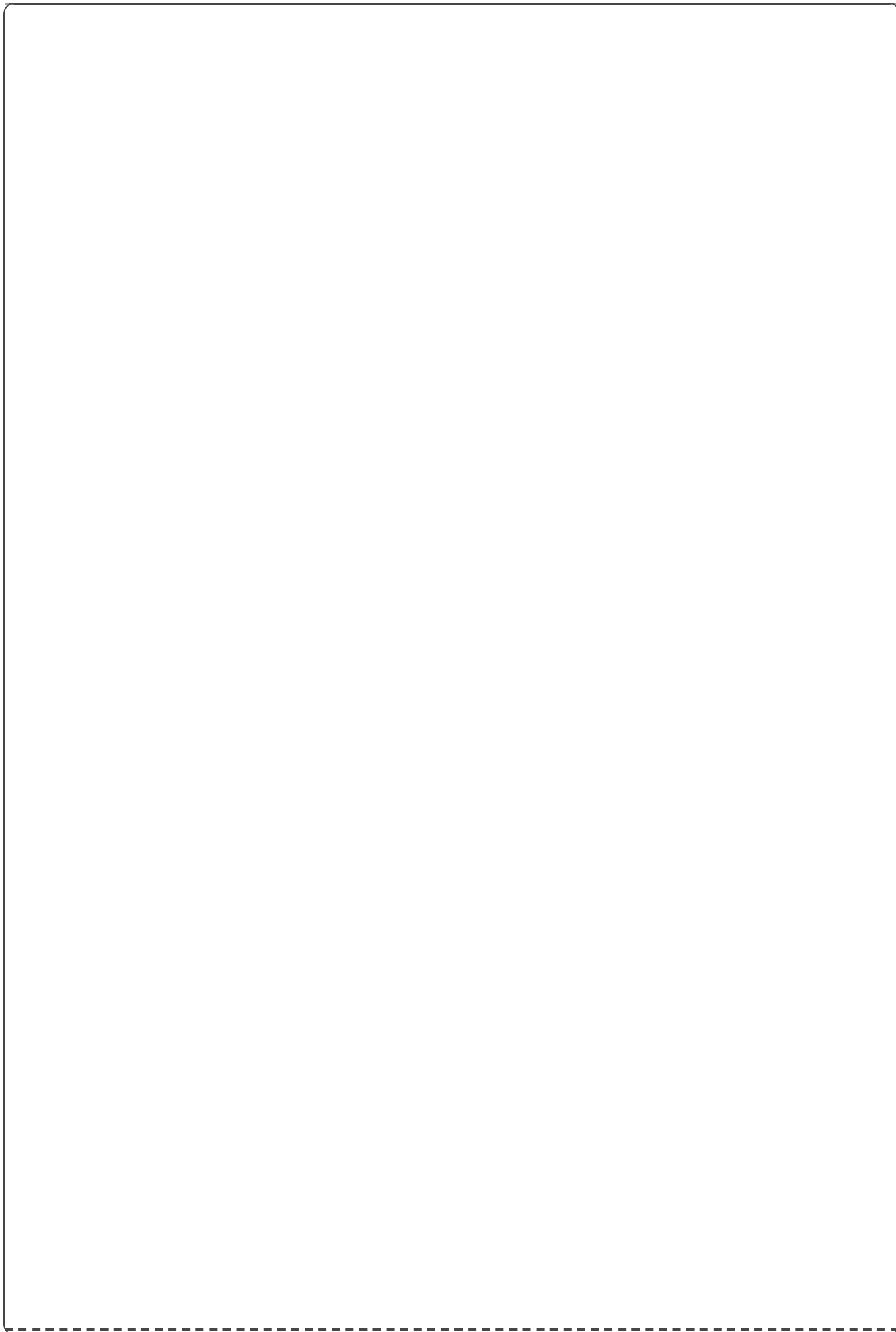
Usar solo la última evaluación y limpieza de texto

```

1 cols_fijas = ["No.", "RUC", "Institucion", "Segmento", "Firma"]
2

```

```
3 ultima_columna = df_ranking.columns[-1]
4
5 # Crear un nuevo DataFrame con las columnas fijas + la última evaluación
6 df_ranking = df_ranking[cols_fijas + [ultima_columna]].copy()
7
8 # --- Función para limpiar y convertir ratings (promedio si hay "|")
9 def limpiar_texto_rating(s):
10     if pd.isna(s):
11         return None
12     s = str(s)
13     s = s.replace("-", "-").replace("-", "-") # Normalizar guiones
14     s = s.upper()
15     s = s.replace("`", "'")
16     s = re.sub(r"\s*\|\s*", " | ", s) # Normalizar separador |
17     s = re.sub(r"\s*-\s*", "-", s) # Normalizar guiones con espacio
18     s = re.sub(r"\s{2,}", " ", s) # Quitar múltiples espacios
19     return s.strip()
20 df_ranking[ultima_columna] = df_ranking[ultima_columna].apply(limpiar_texto_rating)
21 df_ranking
```



Next steps:

[Generate code with df\\_ranking](#)[New interactive sheet](#)

4	2	1090012010001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA L
5	3	0190006247001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA L
6	4	1090056286001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA L
7	5	1790075494001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA L
8	6	1791708040001	CAJA CENTRAL FINANCO
9	7	1390013678001	COOPERATIVA DE AHORRO Y CREDITO 15 DE ABRIL LT
10	8	1790093204001	COOPERATIVA DE AHORRO Y CREDITO 23 DE JULIO LIMITA
11	9	1790567699001	COOPERATIVA DE AHORRO Y CREDITO 29 DE OCTUBRE LT

Cambiar los nombre por números

```

1 rating_scale = {
2     "AAA+": 1, "AAA-": 1,
3     "AA+": 2, "AA": 2, "AA-": 2,
4     "A+": 3, "A": 3, "A-": 3,
5     "BBB+": 4, "BBB": 4, "BBB-": 4,
6     "BB+": 5, "BB": 5, "BB-": 5,
7     "B+": 5, "B": 5, "B-": 5,
8 }
9
10 def convertir_rating(valor, escala):
11     if pd.isna(valor):
12         return None
13
14     # Normalizar texto
15     valor = str(valor).replace("\r", " ").replace("\n", " ").strip()
16

```



```

17 # Convertir cualquier separador (/ | \) en "|"
18 valor = re.sub(r"\s*[ / | \]\s*", "|", valor)
19 valor = re.sub(r"\s+", " ", valor)
20
21 # Caso simple
22 if "|" not in valor:
23     return escala.get(valor, None)
24
25 # Caso con dos ratings
26 partes = [p.strip() for p in valor.split("|")]
27 valores = [escala.get(p, None) for p in partes]
28 valores = [v for v in valores if v is not None]
29
30 if not valores:
31     return None
32
33 return sum(valores) / len(valores)
34
35
36 df_ranking["rating_num"] = df_ranking[ultima_columna].apply(
37     lambda x: convertir_rating(x, rating_scale)
38 )
39
40 # Eliminar filas cuyo rating no existe (None)
41 df_ranking_filtrado = df_ranking[df_ranking["rating_num"].notna()].copy()
42 df_ranking_ordenado = df_ranking_filtrado.sort_values("rating_num").reset_i
43 df_ranking_ordenado
44
45 25 1190075539001 COOPERATIVA DE AHORRO Y CREDITO DE LA PEQUEÑA EMPRE
46
47 26 0590052000001 COOPERATIVA DE AHORRO Y CREDITO DE LA PEQUEÑA EMPRE
48
49 27 1690012606001 COOPERATIVA DE AHORRO Y CREDITO DE LA PEQUEÑA EMPRE

```



				1	6	1791708040001	CAJA CENTRAL FINANCIERA
				Next steps: <a href="#">Generate code with df_ranking_ordenado</a> <a href="#">New interactive sheet</a>			
				2	26	0590052000001	COOPERATIVA DE AHORRO Y CREDITO DE LA PEQUEÑA EMPRE
				3	55	1890003628001	COOPERATIVA DE AHORRO Y CRÉDITO SAN FRANCISCO LT
				4	15	1790325083001	COOPERATIVA DE AHORRO Y CREDITO ANDALUCIA LIMITA
				5	16	1090033456001	COOPERATIVA DE AHORRO Y CREDITO ATUNTAQUI LIMITA
				6	12	1790501469001	COOPERATIVA DE AHORRO Y CREDITO ALIANZA DEL
				7	5	1790075494001	ASOCIACION MUTUALISTA DE AHORRO Y CREDITO PARA L

Indicadores Financieros

				8	29	06900045380001	COOPERATIVA DE AHORRO Y CREDITO RIOBAMBA LT
				1 df_indicadores_financieros			
				9	30	1790979016001	COOPERATIVA DE AHORRO Y CREDITO DE LOS SERVIDORES PI
				10	33	1890037646001	COOPERATIVA DE AHORRO Y CREDITO EL SAGRARIO LT
				11	38	0190155722001	COOPERATIVA DE AHORRO Y CREDITO JARDIN AZUAYO LIM
				12	27	1690012606001	COOPERATIVA DE AHORRO Y CREDITO DE LA PEQUEÑA EMPRI

13	9	1790567699001	COOPERATIVA DE AHORRO Y CRÉDITO 29 DE OCTUBRE LT
14	39	0190115798001	COOPERATIVA DE AHORRO Y CRÉDITO JUVENTUD ECUATORIA
15	3	0190006247001	ASOCIACIÓN MUTUALISTA DE AHORRO Y CRÉDITO PARA LA
16	46	1890001323001	COOPERATIVA DE AHORRO Y CRÉDITO OSCUS LIMITADA
17	56	0290003288001	COOPERATIVA DE AHORRO Y CRÉDITO SAN JOSÉ LIMITADA
18	61	0591711164001	COOPERATIVA DE AHORRO Y CRÉDITO VIRGEN DEL CISO
19	59	0490002669001	COOPERATIVA DE AHORRO Y CRÉDITO TULCÁN LIMITADA
20	35	0691706710001	COOPERATIVA DE AHORRO Y CRÉDITO FERNANDO DAQUIEN
21	34	0190024733001	COOPERATIVA DE AHORRO Y CRÉDITO ERCO LIMITADA

	15 DE ABRIL LTDA	23 DE JULIO LIMITADA	29 DE OCTUBRE LTDA	9 DE OCTUBRE LTDA	ALFONSO JARAMILLO LEON CAJA	1 DEL L1
<b>( PATRIMONIO + RESULTADOS ) / ACTIVOS INMOVILIZADOS</b>	1.320542	2.076457	2.183547	2.788390	1.460808	1
<b>ACTIVOS IMPRODUCTIVOS NETOS / TOTAL ACTIVOS</b>	0.123102	0.091748	0.073416	0.078901	0.094265	0
<b>ACTIVOS PRODUCTIVOS / TOTAL ACTIVOS</b>	0.876898	0.908252	0.926584	0.921099	0.905735	0
<b>ACTIVOS PRODUCTIVOS / PASIVOS CON COSTO</b>	1.027844	1.044289	1.039765	1.080646	1.014125	1
<b>MOROSIDAD DE LA CARTERA DE CREDITO PRODUCTIVO</b>	0.111285	0.020040	0.055395	0.000000	0.080890	0
<b>MOROSIDAD DE LA CARTERA DE CONSUMO</b>	0.070948	0.074157	0.088975	0.036601	0.053926	0
<b>MOROSIDAD DE LA CARTERA DE CREDITO INMOBILIARIO</b>	0.029615	0.101904	0.022915	0.000005	0.000000	0
<b>MOROSIDAD DE LA CARTERA DE MICROCREDITO</b>	0.182316	0.119429	0.196391	0.125266	0.167566	0
<b>MOROSIDAD DE LA CARTERA DE VIVIENDA DE INTERES SOCIAL Y PUBLICO</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0
<b>MOROSIDAD DE LA CARTERA DE CREDITO EDUCATIVO</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0
<b>MOROSIDAD DE LA CARTERA TOTAL</b>	0.107272	0.093425	0.099000	0.096313	0.067223	0

<b>COBERTURA DE LA CARTERA DE CREDITO PRODUCTIVO</b>	1.111456	1.377552	1.000000	0.000000	1.110084	0
<b>COBERTURA DE LA CARTERA DE CREDITO CONSUMO</b>	1.101038	1.133063	1.108068	1.172375	1.022448	1
<b>COBERTURA DE LA CARTERA DE CREDITO INMOBILIARIO</b>	0.631811	0.899309	1.051257	1276.490000	0.000000	563
<b>COBERTURA DE LA CARTERA DE MICROCREDITO</b>	1.007935	1.033881	1.032547	0.934708	1.002167	1
<b>COBERTURA DE LA CARTERA DE VIVIENDA DE INTERES PUBLICO</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0
<b>COBERTURA DE LA CARTERA DE CREDITO EDUCATIVO</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0
<b>COBERTURA DE LA CARTERA PROBLEMÁTICA</b>	1.023424	1.076982	1.080869	0.963912	1.019509	1
<b>GASTOS DE OPERACION ESTIMADOS / TOTAL ACTIVO PROMEDIO</b>	0.055945	0.043137	0.029980	0.038264	0.043149	0
<b>GASTOS DE OPERACION / MARGEN FINANCIERO</b>	1.348749	1.306153	1.051972	1.119722	0.995780	1
<b>GASTOS DE PERSONAL ESTIMADOS / ACTIVO PROMEDIO</b>	0.023108	0.018401	0.012743	0.017336	0.022177	0
<b>RESULTADOS DEL EJERCICIO / PATRIMONIO PROMEDIO</b>	0.005227	0.013795	0.005820	0.006256	0.012872	0
<b>RESULTADOS DEL EJERCICIO / ACTIVO PROMEDIO</b>	0.000585	0.001484	0.000520	0.000738	0.001185	0
<b>CARTERA BRUTA / (DEPOSITOS A LA</b>	0.702486	0.706680	0.856153	0.822414	0.871204	0

<b>VISTA + DEPOSITOS A PLAZO)</b>	0.132400	0.130000	0.000100	0.022414	0.011291	0
<b>MARGEN DE INTERMEDIACIÓN ESTIMADO / PATRIMONIO PROMEDIO</b>	-0.129207	-0.093969	-0.016568	-0.034671	0.001986	-0
<b>MARGEN DE INTERMEDIACIÓN ESTIMADO / ACTIVO PROMEDIO</b>	-0.014466	-0.010111	-0.001481	-0.004091	0.000183	-0
<b>RENDIMIENTO DE LA CARTERA DE CREDITO PRODUCTIVO POR VENCER</b>	0.080513	0.124139	0.083961	0.000000	0.110231	0
<b>RENDIMIENTO DE LA CARTERA DE CREDITO CONSUMO</b>	0.148228	0.139956	0.137721	0.131738	0.149500	0
<b>RENDIMIENTO DE LA CARTERA DE CREDITO INMOBILIARIO POR VENCER</b>	0.098867	0.094859	0.093932	0.090679	0.000000	0
<b>RENDIMIENTO DE LA CARTERA DE MICROCREDITO POR VENCER</b>	0.164794	0.169111	0.163811	0.168581	0.157327	0
<b>RENDIMIENTO DE LA CARTERA DE VIVIENDA DE IINTERES PUBLICO POR VENCER</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0
<b>RENDIMIENTO DE LA CARTERA DE CREDITO EDUCATIVO POR VENCER</b>	0.000000	0.000000	0.000000	0.000000	0.000000	0
<b>CARTERAS DE CRÉDITOS REFINANCIADAS</b>	0.144911	0.162655	0.146202	0.212391	0.197214	0
<b>CARTERAS DE CRÉDITOS REESTRUCTURADAS</b>	0.157736	0.179990	0.137274	0.249426	0.164180	0
<b>CARTERA POR VENCER TOTAL</b>	0.154540	0.161925	0.150601	0.162517	0.151499	0
<b>FONDOS</b>						

DISPONIBLES / TOTAL DEPOSITOS A CORTO PLAZO	2014	2015	2016	2017	2018	2019	2020	2021	2022
Hay que dividir los datos para el 2022, porque en el documento aparece con el 0	0.1242	0.1331	0.2784	0.1646	0.2035	0.2157	0.1748	0.1574	0.17748

CARTERA

```
1 # Copia del DataFrame original
2 df_fixed = df_indicadores_financieros.copy()
3
4 # 1. Reemplazar % y convertir todo a número (coercing evita errores)
5 df_fixed = df_fixed.replace('%', '', regex=True)
6 df_fixed = df_fixed.apply(pd.to_numeric, errors='coerce')
7
8 # 2. Convertir porcentajes (100% -> 1.0)
9 df_fixed = df_fixed / 100
10
11 # 3. Mostrar primeras filas
12 df_fixed.head()
```

FI = 1 + (ACTIVOS IMPRODUCTIVOS / ACTIVOS TOTALES)	1.123102	1.091748	1.073416	1.078901	1.094265	1.094265	1.094265
	15 DE ABRIL	23 DE JULIO	29 DE OCTUBRE	9 DE OCTUBRE	ALFONSO JARAMILLO	ALIANZA DEL VALLE	ALIANZA DEL VALLE
INDICE DE CAPITALIZACION NETO: FK / FI	LTDA	LIMITADA	LTDA	LTDA	LEON CAJA	LIMITADA	LIMITADA
	0.082552	0.088861	0.078577	0.109300	0.078634	0.078634	0.078634
( PATRIMONIO + RESULTADOS ) / ACTIVOS INMOVILIZADOS	0.013205	0.020765	0.021835	0.027884	0.014608	0.016484	0.016484
ACTIVOS IMPRODUCTIVOS NETOS / TOTAL ACTIVOS	0.001231	0.000917	0.000734	0.000789	0.000943	0.000945	0.000945
ACTIVOS PRODUCTIVOS / TOTAL ACTIVOS	0.008769	0.009083	0.009266	0.009211	0.009057	0.009055	0.009055
ACTIVOS PRODUCTIVOS / PASIVOS CON COSTO	0.010278	0.010443	0.010398	0.010806	0.010141	0.010461	0.010461
MOROSIDAD DE LA CARTERA DE CREDITO PRODUCTIVO	0.001113	0.000200	0.000554	0.000000	0.000809	0.000000	0.000000

5 rows × 44 columns



```
1 !pip install unidecode
2
```

Requirement already satisfied: unidecode in /usr/local/lib/python3.12/dist-packa

```
1 import unidecode
2
3
4 def normalizar_texto(valor):
5     if not isinstance(valor, str):
6         return valor # no modificar números
7     valor = valor.lower()
8     valor = unidecode.unidecode(valor)
9     valor = re.sub(r'^a-z0-9+', '_', valor)
10    valor = re.sub(r'_+', '_', valor)
11    valor = valor.strip('_')
12    return valor
13
14 df_indicadores_financieros.index = df_indicadores_financieros.index.map(nor
15
16 df_indicadores_financieros.columns = df_indicadores_financieros.columns.map
17
18 print("Normalización completada para índices y columnas de df_indicadores_f
19 display(df_indicadores_financieros.head())
```

Normalización completada para índices y columnas de df\_indicadores\_financieros.

15\_de\_abril ltda 23\_de\_julio limit

<b>patrimonio_resultados_activos_inmovilizados</b>	1.320542	2.076
<b>activos_improductivos_netos_total_activos</b>	0.123102	0.091
<b>activos_productivos_total_activos</b>	0.876898	0.908
<b>activos_productivos_pasivos_con_costo</b>	1.027844	1.044
<b>morosidad_de_la_cartera_de_credito_productivo</b>	0.111285	0.020

5 rows × 44 columns

```
1 df_indicadores_financieros.index
```

```
Index(['patrimonio_resultados_activos_inmovilizados',
      'activos_improductivos_netos_total_activos',
      'activos_productivos_total_activos',
      'activos_productivos_pasivos_con_costo',
      'morosidad_de_la_cartera_de_credito_productivo',
      'morosidad_de_la_cartera_de_consumo',
      'morosidad_de_la_cartera_de_credito_inmobiliario',
      'morosidad_de_la_cartera_de_microcredito',
```

```
'morosidad_de_la_cartera_de_vivienda_de_interes_social_y_publico',
'morosidad_de_la_cartera_de_credito_educativo',
'morosidad_de_la_cartera_total',
'cobertura_de_la_cartera_de_credito_productivo',
'cobertura_de_la_cartera_de_credito_consumo',
'cobertura_de_la_cartera_de_credito_inmobiliario',
'cobertura_de_la_cartera_de_microcredito',
'cobertura_de_la_cartera_de_vivienda_de_iinteres_publico',
'cobertura_de_la_cartera_de_credito_educativo',
'cobertura_de_la_cartera_problematika',
'gastos_de_operacion_estimados_total_activo_promedio',
'gastos_de_operacion_margen_financiero',
'gastos_de_personal_estimados_activo_promedio',
'resultados_del_ejercicio_patrimonio_promedio',
'resultados_del_ejercicio_activo_promedio',
'cartera_bruta_depositos_a_la_vista_depositos_a_plazo',
'margen_de_intermediacion_estimado_patrimonio_promedio',
'margen_de_intermediacion_estimado_activo_promedio',
'rendimiento_de_la_cartera_de_credito_productivo_por_vencer',
'rendimiento_de_la_cartera_de_credito_consumo',
'rendimiento_de_la_cartera_de_credito_inmobiliario_por_vencer',
'rendimiento_de_la_cartera_de_microcredito_por_vencer',
'rendimiento_de_la_cartera_de_vivienda_de_iinteres_publico_por_vencer',
'rendimiento_de_la_cartera_de_credito_educativo_por_vencer',
'carteras_de_creditos_refinanciadas',
'carteras_de_creditos_reestructuradas', 'cartera_por_vencer_total',
'fondos_disponibles_total_depositos_a_corto_plazo',
'cartera_improductiva_descubierta_patrimonio_resultados',
'cartera_improductiva_patrimonio_dic',
'fk_patrimonio_resultados_ingresos_extraordinarios_activos_totales',
'fi_1_activos_improductivos_activos_totales',
'indice_de_capitalizacion_neto_fk_fi'],
dtype='object', name='')
```

## ✓ Escalado con Standard Scaler

```
1 df_indicadores_financieros
```



15\_de\_abril\_

patrimonio_resultados_activos_inmovilizados	1.32
activos_improductivos_netos_total_activos	0.12
activos_productivos_total_activos	0.87
activos_productivos_pasivos_con_costo	1.02
morosidad_de_la_cartera_de_credito_productivo	0.11
morosidad_de_la_cartera_de_consumo	0.07
morosidad_de_la_cartera_de_credito_inmobiliario	0.02
morosidad_de_la_cartera_de_microcredito	0.18
morosidad_de_la_cartera_de_vivienda_de_interes_social_y_publico	0.00
morosidad_de_la_cartera_de_credito_educativo	0.00
morosidad_de_la_cartera_total	0.10
cobertura_de_la_cartera_de_credito_productivo	1.11

```

1 from sklearn.preprocessing import StandardScaler
2
3 X = df_indicadores_financieros.copy()
4
5 # Convertir todo a numérico por seguridad
6 X = X.apply(pd.to_numeric, errors='coerce')
7
8 # Llenar valores faltantes
9 X = X.fillna(0)
10
11 # Escalar
12 scaler = StandardScaler()
13 X_scaled = scaler.fit_transform(X)
14
15 df_indicadores_financieros_scaled = pd.DataFrame(
16     X_scaled,
17     index=X.index,
18     columns=X.columns
19 )
20
21 df_indicadores_financieros_scaled

```

margen_de_intermediacion_estimado_activo_promedio	-0.01
rendimiento_de_la_cartera_de_credito_productivo_por_vencer	0.08
rendimiento_de_la_cartera_de_credito_consumo	0.14
rendimiento_de_la_cartera_de_credito_inmobiliario_por_vencer	0.09
rendimiento_de_la_cartera_de_microcredito_por_vencer	0.16

rendimiento_de_la_cartera_de_vivienda_de_iinteres_publico_por_vencer	0.00
rendimiento_de_la_cartera_de_credito_educativo_por_vencer	0.00
carteras_de_creditos_refinanciadas	0.14
carteras_de_creditos_reestructuradas	0.15
cartera_por_vencer_total	0.15
fondos_disponibles_total_depositos_a_corto_plazo	0.24
cartera_improductiva_descubierta_patrimonio_resultados	0.00
cartera_improductiva_patrimonio_dic	0.65
fk_patrimonio_resultados_ingresos_extraordinarios_activos_totales	0.09
fi_1_activos_improductivos_activos_totales	1.12
indice_de_capitalizacion_neto_fk-fi	0.08

41 rows × 44 columns

15\_de\_abril\_

patrimonio_resultados_activos_inmovilizados	2.19
activos_improductivos_netos_total_activos	-0.48
activos_productivos_total_activos	1.20
activos_productivos_pasivos_con_costo	1.53
morosidad_de_la_cartera_de_credito_productivo	-0.51
morosidad_de_la_cartera_de_consumo	-0.60
morosidad_de_la_cartera_de_credito_inmobiliario	-0.69
morosidad_de_la_cartera_de_microcredito	-0.35
morosidad_de_la_cartera_de_vivienda_de_interes_social_y_publico	-0.76
morosidad_de_la_cartera_de_credito_educativo	-0.76
morosidad_de_la_cartera_total	-0.52
cobertura_de_la_cartera_de_credito_productivo	1.72

>	<b>Análisis exploratorio (EDA)</b>	cobertura_de_la_cartera_de_credito_consumo	1.70
		cobertura_de_la_cartera_de_credito_inmobiliario	0.65
	↳ 11 cells hidden	cobertura_de_la_cartera_de_microcredito	1.49
		cobertura_de_la_cartera_de_vivienda_de_iinteres_publico	-0.76
>	<b>TSNE</b>	cobertura_de_la_cartera_de_credito_educativo	-0.76
		cobertura_de_la_cartera_problematica	1.52
	↳ 5 cells hidden	gastos_de_operacion_estimados_total_activo_promedio	-0.63
		gastos_de_operacion_margen_financiero	2.25
✓	<b>Modelado</b>	gastos_de_personal_estimados_activo_promedio	-0.71
		resultados_del_ejercicio_patrimonio_promedio	-0.75
>	<b>KMeans</b>	resultados_del_ejercicio_activo_promedio	-0.76
		cartera_bruta_depositos_a_la_vista_depositos_a_plazo	1.01
	↳ 13 cells hidden	margen_de_intermediacion_estimado_patrimonio_promedio	-1.05
		margen_de_intermediacion_estimado_activo_promedio	-0.79
>	<b>Gaussian Mixture Model</b>	rendimiento_de_la_cartera_de_credito_productivo_por_vencer	-0.58
		rendimiento_de_la_cartera_de_credito_consumo	-0.43
	↳ 4 cells hidden	rendimiento_de_la_cartera_de_credito_inmobiliario_por_vencer	-0.54
		rendimiento_de_la_cartera_de_microcredito_por_vencer	-0.39

## Spectral Clustering

rendimiento\_de\_la\_cartera\_de\_vivienda\_de\_iinteres\_publico\_por\_vencer

-0.76

rendimiento\_de\_la\_cartera\_de\_credito\_educativo\_por\_vencer

0.76

```

1 from sklearn.cluster import SpectralClustering
2 from sklearn import metrics
3
4 k_range = range(2, 11)
5 silhouette_scores = []
6
7 print("Calculando Silueta para Spectral Clustering...")
8 for k in k_range:
9     # Usamos los mismos parámetros que el ajuste anterior
10    spectral_clustering = SpectralClustering(
11        n_clusters=k,
12        assign_labels='kmeans',
13        affinity='rbf',
14        gamma=1.0, # Ajusta esto si es necesario
15        random_state=42,
16        n_init=10
17    )
18    # Ajustar y predecir
19    labels = spectral_clustering.fit_predict(X)
20
21    # Calcular la Silueta (solo si k > 1)
22    score = metrics.silhouette_score(X, labels)
23    silhouette_scores.append(score)
24
25 # Encontrar el k con el score de Silueta más alto
26 optimal_k_spectral = k_range[np.argmax(silhouette_scores)]
27 print(f"El k óptimo para Spectral Clustering (máxima Silhouette) es: {optimal_k_spectral}")
28
29 # Graficar la Silueta
30 plt.figure(figsize=(10, 5))
31 plt.plot(k_range, silhouette_scores, marker='o', linestyle='--')
32 plt.title("Coeficiente de Silueta para Spectral Clustering")
33 plt.xlabel("Número de Clusters (k)")
34 plt.ylabel("Puntuación de Silueta")
35 plt.xticks(k_range)
36 plt.grid(True)
37 plt.show()

```

Calculando Silueta para Spectral Clustering...

El k óptimo para Spectral Clustering (máxima Silhouette) es: 4



Como en los anteriores se ligio 5, para que sea más comparable, en este igual se trabajara con 5 clusters porque tenemos 5 divisiones de cooperativas. Aunque aquí nos indique que lo más óptimo es 4.

```

1 n_neighbors_candidates = [5, 8, 10, 12, 15, 20]
2 from sklearn.metrics import adjusted_rand_score, silhouette_score
3
4 for nn in n_neighbors_candidates:
5     sc = SpectralClustering(
6         n_clusters=k,
7         affinity="nearest_neighbors",
8         n_neighbors=nn,
9         assign_labels="kmeans",
10        random_state=42
11    )
12    labels = sc.fit_predict(X_tsne)
13    ari = adjusted_rand_score(labels_K, labels)
14    sil = silhouette_score(X_tsne, labels)
15    print(f"n_neighbors={nn} ARI={ari:.4f} Silhouette={sil:.4f}")

```



```

/usr/local/lib/python3.12/dist-packages/sklearn/manifold/_spectral_embedding.py:
warnings.warn(
n_neighbors=5  ARI=0.5211  Silhouette=0.5131
n_neighbors=8  ARI=0.5592  Silhouette=0.4130
n_neighbors=10  ARI=0.4860  Silhouette=0.3283
n_neighbors=12  ARI=0.4943  Silhouette=0.3917
n_neighbors=15  ARI=0.4640  Silhouette=0.2337
n_neighbors=20  ARI=0.4185  Silhouette=0.1795

```

Como el ARI más alto y el Silhouette mas alto es cuando hay 8 neighbors, se elegiran 8.

```

1 from sklearn.cluster import SpectralClustering
2
3 # Número de clusters
4 k_spectral = 5
5
6 # Definir Spectral Clustering sobre TSNE
7 spectral_tsne = SpectralClustering(
8     n_clusters=k_spectral,
9     affinity='nearest_neighbors',
10    n_neighbors= 8,
11    assign_labels='kmeans',
12    random_state=42,
13    n_init=10
14 )
15
16 labels_spectral= spectral_tsne.fit_predict(X_tsne)
17
18 # Crear dataframe copiando indicadores y añadiendo clusters spectral sobre
19 df_clusters_spectral = df_indicadores_financieros_scaled_trasp.copy()
20 df_clusters_spectral["cluster_spectral_tsne"] = labels_spectral
21
22 # Ver distribución
23 print(df_clusters_spectral["cluster_spectral_tsne"].value_counts())

```

```

cluster_spectral_tsne
0     11
3     10
2     10
4      7
1      6
Name: count, dtype: int64

```

```

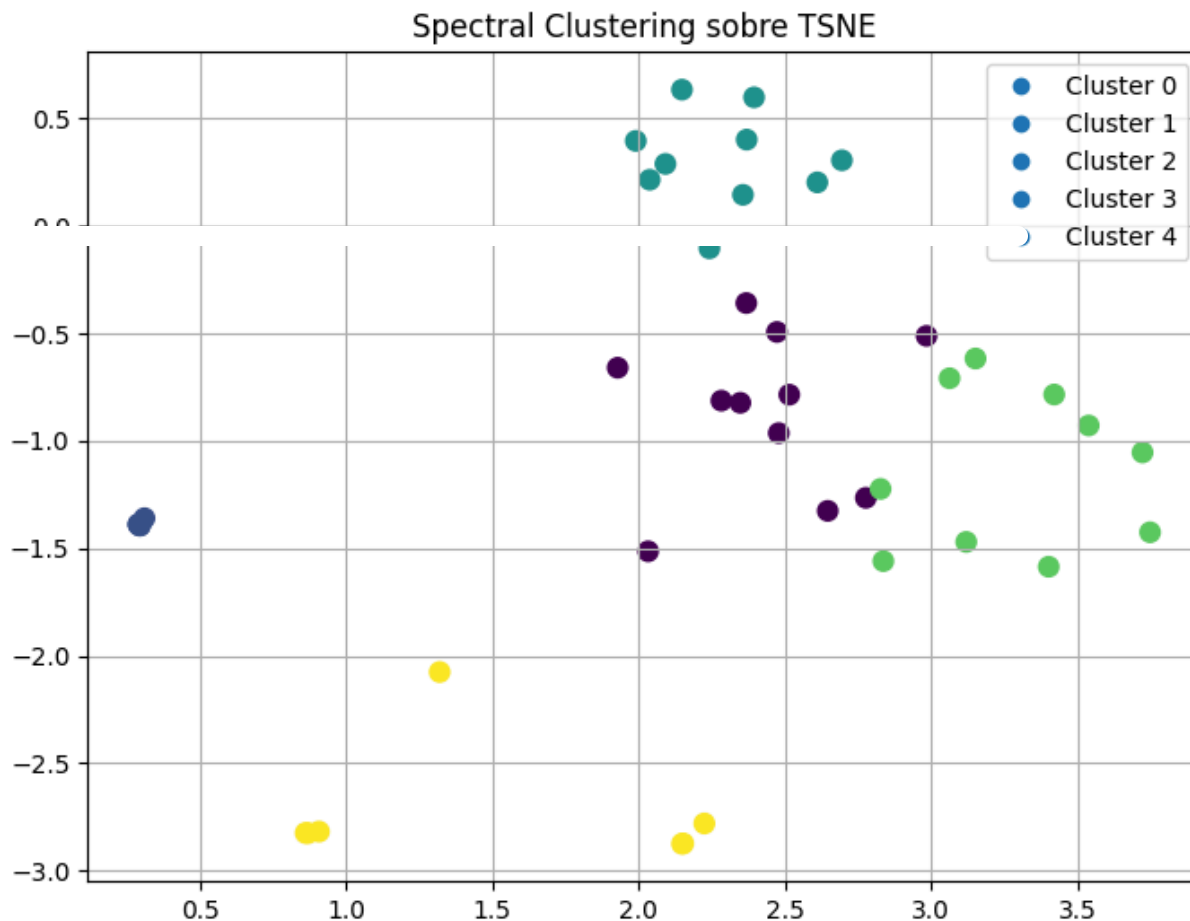
1
2 plt.figure(figsize=(8,6))
3 scatter = plt.scatter(X_tsne[:,0], X_tsne[:,1], c=labels_spectral, s=60)
4 plt.title("Spectral Clustering sobre TSNE")
5 plt.grid(True)
6
7 # Crear leyenda con los nombres de los clusters

```

```

8 handles = []
9 unique_labels = np.unique(labels_spectral)
10 for lab in unique_labels:
11     handles.append(plt.Line2D([], [], marker='o', linestyle='', label=f'Clu
12
13 plt.legend(handles=handles, loc='best')
14 plt.show()

```

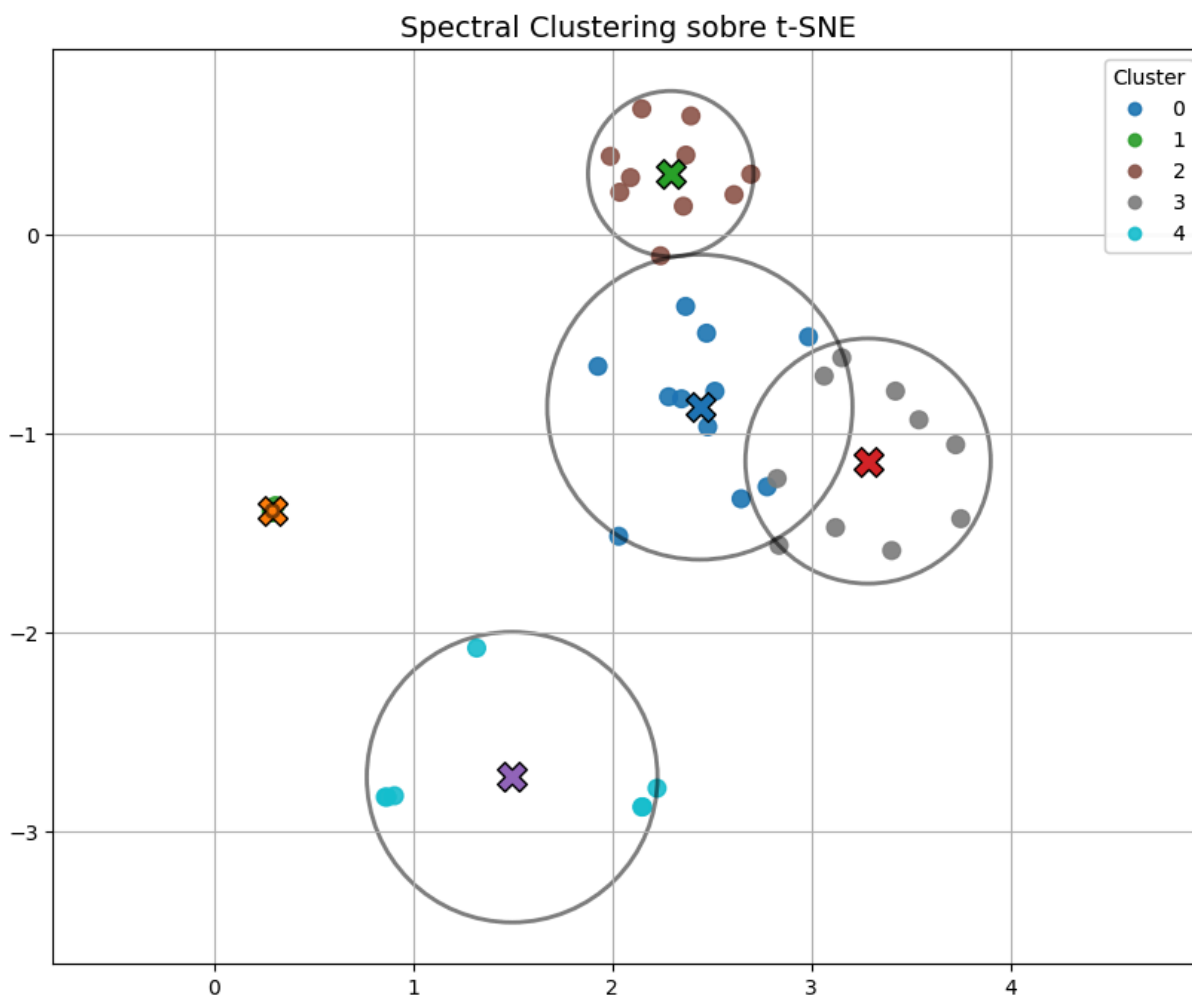


```

1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 plt.figure(figsize=(10, 8))
5
6 # Scatter con mejor colormap
7 scatter = plt.scatter(
8     X_tsne[:, 0],
9     X_tsne[:, 1],
10    c=labels_spectral,
11    cmap="tab10",
12    s=70,
13    alpha=0.9
14 )

```

```
15
16 # Obtener centroides y dibujar círculos reales
17 for cluster in np.unique(labels_spectral):
18
19     pts = X_tsne[labels_spectral == cluster]
20     centroid = pts.mean(axis=0)
21
22     # Dibujar centroide
23     plt.scatter(
24         centroid[0],
25         centroid[1],
26         marker="X",
27         s=200,
28         edgecolors="black"
29     )
30
31     # Radio más estable: distancia máxima (no promedio)
32     dx = pts[:, 0] - centroid[0]
33     dy = pts[:, 1] - centroid[1]
34     radius = np.max(np.sqrt(dx**2 + dy**2))
35
36     circle = plt.Circle(
37         (centroid[0], centroid[1]),
38         radius,
39         edgecolor="black",
40         fill=False,
41         linewidth=2,
42         alpha=0.5
43     )
44
45     plt.gca().add_patch(circle)
46
47 # Título, grilla y formato
48 plt.title("Spectral Clustering sobre t-SNE", fontsize=14)
49 plt.grid(True)
50 plt.gca().set_aspect("equal", adjustable="datalim")
51
52 # Leyenda automática
53 legend1 = plt.legend(*scatter.legend_elements(), title="Cluster")
54 plt.gca().add_artist(legend1)
55
56 plt.show()
57
```



## ✓ Evaluación y validación

## ✓ Comparación de Resultados

Silhouette Score, mide cohesión (qué tan juntos están) y separación (qué tan lejos están de otros clusters). Aquí mientras más cercana a 1 mejor.

```

1 from sklearn.metrics import silhouette_score
2
3 sil_km = silhouette_score(X, labels_K)
4 sil_gmm = silhouette_score(X, labels_gmm)
5 sil_spectral = silhouette_score(X, labels_spectral)
6
7 print("Silhouette K-Means:", sil_km)
8 print("Silhouette GMM:", sil_gmm)
9 print("Silhouette Spectral:", sil_spectral)

```

```

Silhouette K-Means: 0.3724289016235934
Silhouette GMM: 0.4035526757905385
Silhouette Spectral: 0.37524506219056913

```

Con Silhouette se sacaron resultados muy pobres, mientras que los mejores fueron con KMeans, sin decir que fueron escalantes

El metodo de Clinski - Harabasz mide la separación entre clusters (between-cluster dispersion), y cohesión dentro del cluster (within-cluster dispersion). Mientras más alto es mejor.

```

1 from sklearn.metrics import calinski_harabasz_score
2
3 ch_km = calinski_harabasz_score(X, labels_K)
4 ch_gmm = calinski_harabasz_score(X, labels_gmm)
5 ch_spectral = calinski_harabasz_score(X, labels_spectral)
6
7 print("CH K-Means:", ch_km)
8 print("CH GMM:", ch_gmm)
9 print("CH Spectral:", ch_spectral)

```

```

CH K-Means: 24.936789828579904
CH GMM: 23.726184628838197
CH Spectral: 25.00650637750751

```

Con esta métrica también se coincide que KMeans es la mejor y Spectral por mucho la peor.

## ✓ Mutual Information

Mide qué tan informativo es cada indicador financiero para explicar el rating oficial de riesgo. Si saber el valor de un indicador reduce la incertidumbre sobre el rating entonces tiene Alta Información Mutua. Si no aporta nada su Información Mutua es cercana a cero

El código calcula la Información Mutua entre cada indicador financiero y el rating de riesgo oficial, permitiendo identificar qué variables aportan más información útil para explicar o predecir el rating.

## ✓ Adjusted Rand Index (ARI)

1 Start coding or generate with AI.

## ✓ Conclusiones

## ✓ Similitudes y discrepancias entre clusters y ratings.

Definitivamente el que hizo una separación más pareja fue el KMeans, divide claramente en 5 grupos compactos. A simple vista no maneja bien los grupos internos ni alargados, pero es el mejor de los 3 modelos.

Con el GMM hay un cluster gigante (celeste) que absorbe muchos puntos. Algunos clusters están parcialmente solapados. Indica que hay cooperativas con comportamientos financieros muy parecidos

En el Spectral teóricamente se detectan estructuras no lineales, parece que separa otros clusters. Si se hubiera elegido  $k = 4$ , que es lo que en la mayoría de modelos sugería, hubiera habido un rendimiento un poco mejor, pero se hubieran tenido que dividir las cooperativas de otra manera

Hipótesis de los patrones financieros.

## ✓ Semi Supervised Learning

## ✓ Alineación de nombres + Matching Ratings–Clusters

Objetivo: Unificar los nombres de cooperativas para poder comparar clusters financieros vs. ratings oficiales.

- Se crearon claves normalizadas (Cooperativa\_Key, Institucion\_Key)

- Se ejecutó RapidFuzz con partial\_ratio (correcto para nombres largos e inconsistentes)
- Se construyó un mapeo confiable ( $\geq 90$  de score)
- Se generó df\_validacion con 42 coincidencias exactas
- Ya tenemos la tabla emparejada con: Cooperativa — Cluster — rating\_num



```
1 !pip install rapidfuzz
```

```
Requirement already satisfied: rapidfuzz in /usr/local/lib/python3.12/dist-packa
```

### Construir df\_cluster\_labels

```
1 # =====
2 # RECONSTRUIR df_cluster_labels (como en tu parte 1)
3 # =====
4
5 # df_clusters_km tiene el índice = cooperativa normalizada
6 # y la columna cluster_km
7
8 df_cluster_labels = df_clusters_km.reset_index().rename(columns={
9     "index": "Cooperativa",
10    "cluster_km": "Cluster"
11 })
12
13 # Dejamos solo las columnas necesarias para el matching
14 df_cluster_labels = df_cluster_labels[["Cooperativa", "Cluster"]]
15
16 print("✓ df_cluster_labels creado correctamente")
17 display(df_cluster_labels.head())
18
```

✓ df\_cluster\_labels creado correctamente

	Cooperativa	Cluster	
0	15_de_abril ltda	2	
1	23_de_julio limitada	3	
2	29_de_octubre ltda	3	
3	9_de_octubre ltda	1	
4	alfonso_jaramillo_leon_caja	2	

```
1 # =====
2 # BLOQUE 1 – Normalización de nombres + Matching
```

```

3 # =====
4
5 import unicodedata
6 import re
7 from rapidfuzz import process, fuzz
8
9 # Función para normalizar nombres (quita tildes, espacios, mayúsculas, etc.
10 def normalizar(texto):
11     texto = texto.lower()
12     texto = unicodedata.normalize("NFKD", texto).encode("ascii", "ignore").
13     texto = re.sub(r'^a-z0-9+', '_', texto)
14     texto = texto.strip('_')
15     return texto
16
17 # -----
18 # 1) Crear claves normalizadas en ambos DFs
19 # -----
20 df_cluster_labels["Cooperativa_Key"] = df_cluster_labels["Cooperativa"].app
21 df_ranking_ordenado["Institucion_Key"] = df_ranking_ordenado["Institucion"]
22
23 print("✓ Claves normalizadas creadas correctamente.")
24 print(df_cluster_labels.head())
25 print(df_ranking_ordenado.head())
26
27 # -----
28 # 2) Matching automático con RapidFuzz
29 # -----
30 cluster_keys = df_cluster_labels["Cooperativa_Key"].tolist()
31 rating_keys = df_ranking_ordenado["Institucion_Key"].tolist()
32
33 matches = []
34
35 print("\n🔍 Buscando match automático entre cooperativas...")
36 for r in rating_keys:
37     best_match, score, _ = process.extractOne(
38         r,
39         cluster_keys,
40         scorer=fuzz.partial_ratio
41     )
42
43     matches.append({
44         "rating_key": r,
45         "best_cluster_key": best_match,
46         "score": score
47     })
48
49 df_matches = pd.DataFrame(matches).sort_values(by="score", ascending=False)
50
51 print("\n✓ 20 mejores coincidencias:")
52 display(df_matches.head(20))
53

```



```
54 # -----
55 # 3) Filtramos coincidencias confiables (score ≥ 90)
56 # -----
57 df_good = df_matches[df_matches["score"] >= 90].copy()
58
59 map_dict = dict(zip(df_good["rating_key"], df_good["best_cluster_key"]))
60
61 print("\n✓ Diccionario final de mapeo (rating → cluster):")
62 map_dict
63
64 # -----
65 # 4) Aplicar el mapeo al DF de ratings
66 # -----
67 df_ranking_ordenado["Cooperativa_Key_Traducida"] = \
68     df_ranking_ordenado["Institucion_Key"].map(map_dict)
69
70 # -----
71 # 5) Generar df_validacion (coincidencias exactas)
72 # -----
73 df_validacion = pd.merge(
74     df_cluster_labels,
75     df_ranking_ordenado,
76     left_on="Cooperativa_Key",
77     right_on="Cooperativa_Key_Traducida",
78     how="inner"
79 )
80
81 print("\n✓ Filas emparejadas:", len(df_validacion))
82 display(df_validacion[["Cooperativa", "Cluster", "rating_num"]])
83
```



```
✓ Claves normalizadas creadas correctamente.

Cooperativa Cluster Cooperativa_Key
0 15_de_abril ltda 2 15_de_abril ltda
1 23_de_julio limitada 3 23_de_julio limitada
2 29_de_octubre ltda 3 29_de_octubre ltda
3 9_de_octubre ltda 1 9_de_octubre ltda
4 alfonso_jaramillo_leon_caja 2 alfonso_jaramillo_leon_caja

No. RUC Institucion \
0 1 1768168480001 CORPORACIÓN NACIONAL DE FINANZAS \r POPULARES Y SOL...
1 6 1791708040001 CAJA CENTRAL FINANCOOP
2 26 0590052000001 COOPERATIVA DE AHORRO Y CREDITO DE LA \r PEQUEÑA EMPRE...
3 55 1890003628001 COOPERATIVA DE AHORRO Y CRÉDITOS AN \r FRANCISCO LTDA.
4 15 1790325083001 COOPERATIVA DE AHORRO Y CREDITO \r ANDALUCIA LIMITADA
```

```
Segmento Firma Jun2025 \
0 1 PACIFIC CREDIT RATING S.A. AAA
1 1 PACIFIC CREDIT RATING S.A. AAA-
2 1 PACIFIC CREDIT RATING S.A. AAA-
3 1 PACIFIC CREDIT RATING S.A. /\r SUMMA RATINGS S.A. AAA-/
4 1 CLASS INTERNATIONAL RATING /\r SUMMA RATINGS S.A. AA+
```

```
rating_num Institucion_Key
0 1.0 corporacionnacionaldefinanzas_populares_y_soli...
1 1.0 caja_central_financoop
2 1.0 cooperativadeahorroycredito dela_pequena_empres...
3 1.0 cooperativadeahorroycreditosan_francisco ltda
4 2.0 cooperativadeahorroycredito_andalucia_limitada
```

🔍 Buscando match automático entre cooperativas...


✓ 20 mejores coincidencias:

	rating_key	best_cluster_key	sco
3	cooperativadeahorroycreditosan_francisco ltda	san_francisco ltda	100
6	cooperativa_de_ahorro_y_credito_alianza_del_valle...	alianza_del_valle limitada	100
5	cooperativadeahorroycredito_atuntaqui limitada	atuntaqui limitada	100
4	cooperativadeahorroycredito_andalucia limitada	andalucia limitada	100
10	cooperativadeahorroycreditoel_sagrario ltda	el_sagrario ltda	100
11	cooperativadeahorroycreditojardin_azuayo limitada	jardin_azuayo limitada	100
8	cooperativadeahorroycredito_riobamba ltda	riobamba ltda	100
30	cooperativa_de_ahorro_y_credito_ambato ltda	ambato ltda	100
24	cooperativa_de_ahorro_y_credito_mushuc_runa ltda	mushuc_runa ltda	100
22	cooperativadeahorroycreditopolicia_nacional_li...	policia_nacional limitada	100
21	cooperativadeahorroycreditoerco limitada	erco limitada	100
20	cooperativadeahorroycredito_fernando_daquilema...	fernando_daquilema limitada	100
19	cooperativadeahorro_ycredito_tulcan limitada	tulcan limitada	100

18	cooperativadeahorroycreditovirgen_del_cisne	virgen_del_cisne	100
17	cooperativadeahorroycreditosan_jose_limitada	san_jose_limitada	100
16	cooperativadeahorroycreditooscus_limitada	oscus_limitada	100
47	cooperativadeahorroycreditopablo_munoz_vega_li...	pablo_munoz_vega_limitada	100
41	cooperativa_de_ahorro_y_credito_once_de_junio_...	once_de_junio ltda	100
42	cooperativadeahorroycreditosanta_rosa_limitada	santa_rosa_limitada	100
43	cooperativadeahorroycreditola_merced_limitada	la_merced_limitada	100

✓ Diccionario final de mapeo (rating → cluster):

✓ Filas emparejadas: 42

	Cooperativa	Cluster	rating_num	
0	15_de_abril ltda	2	5.0	
1	23_de_julio_limitada	3	3.0	
2	29_de_octubre ltda	3	2.0	
3	9_de_octubre ltda	1	3.0	
4	alfonso_jaramillo_leon_caja	2	4.0	
5	alianza_del_valle_limitada	1	2.0	
6	ambato ltda	0	3.0	
7	andalucia_limitada	2	2.0	
8	atuntaqui_limitada	3	2.0	
9	calceta ltda	2	4.0	
10	chibuleo_limitada	2	4.0	
11	chone ltda	0	4.0	
12	comercio ltda	4	3.0	
13	cooprogreso_limitada	4	4.0	
14	de_la_pequena_empresa_biblian_limitada	3	3.0	
15	de_la_pequena_empresa_cacpe_loja ltda	1	3.0	
16	de_la_pequena_empresa_de_cotopaxi_limitada	4	1.0	
17	de_la_pequena_empresa_de_pastaza_limitada	0	2.0	
18	de_la_pequena_empresa_gualaquiza	1	4.0	
19	de_los_servidores_publicos_del_ministerio_de_e...	0	2.0	
20	el_sagrario ltda	0	2.0	

## ✓ CÁLCULO DE MÉTRICAS DE SIMILITUD ENTRE CLUSTERS Y RATINGS

22	fernando_daquilema_limitada	3	3.0
(Adjusted Rand Index y Normalized Mutual Information)			
23	indigena_sac Ltda	2	4.0

```

1 # =====
2 # BLOQUE 2 – Métricas de similitud Clustering vs Rating
3 # =====
4
5 from sklearn.metrics import adjusted_rand_score, normalized_mutual_info_score
6
7 # Extraer vectores alineados
8 labels_pred = df_validacion["Cluster"].astype(int)
9 labels_true = df_validacion["rating_num"].astype(int)
10
11 # Calcular métricas
12 ari = adjusted_rand_score(labels_true, labels_pred)
13 nmi = normalized_mutual_info_score(labels_true, labels_pred)
14
15 print("🔵 Evaluación Clustering vs Ratings Oficiales")
16 print(f"Adjusted Rand Index (ARI): {ari:.4f}")
17 print(f"Normalized Mutual Information (NMI): {nmi:.4f}")
18

```

36	san_francisco Ltda	4	1.0
Adjusted Rand Index (ARI): 0.1566			
37	san_jose_limitada	0	2.0

- 38 santa\_rosa\_limitada 2 4.0
- ARI = 0.1566 → los clusters capturan una fracción pequeña, pero no despreciable, de la estructura real de los ratings. La alineación existe, pero es débil.
- 39 tulcan\_limitada 0 2.0
- NMI = 0.2703 → hay una cantidad moderada de información compartida entre clusters y ratings, indicando que los agrupamientos financieros reflejan parcialmente los patrones que las agencias calificadoras consideran.
- 40 vicentina\_manuel\_esteban\_godoy\_arrieta\_limitada 1 3.0
- 41 virgin\_quecunche 0 2.0

## ✓ APRENDIZAJE SEMISUPERVISADO

“Usaremos el mismo conjunto de variables financieras. Los labels serán los ratings oficiales...”

Ahora empezamos a construir el dataset final (features + label) que usarán TODOS los modelos: baseline, self-training, label spreading.

### ✓ 2.1 Construcción del dataset etiquetado (X\_final, y\_final)

Instruccion: "Reutilizar escalado, imputación y selección de variables definidos en la PARTE 1."

"Los labels serán los ratings oficiales de cada cooperativa."

```

1 # =====
2 # 2.1 DATASET FINAL PARA APRENDIZAJE SEMISUPERVISADO
3 # =====
4
5 import pandas as pd
6
7 # X: mismas variables financieras escaladas de la PARTE 1
8 # (df_indicadores_financieros_scaled_trasp: filas = cooperativas, columnas
9 X_full = df_indicadores_financieros_scaled_trasp.copy()
10
11 # y: ratings oficiales alineados (desde df_validacion)
12 df_labels = df_validacion[["Cooperativa", "rating_num"]].drop_duplicates()
13
14 # Unimos: solo cooperativas que tienen indicadores + rating
15 df_final = X_full.merge(
16     df_labels,
17     left_index=True,          # índice de X_full = nombre de cooperativa
18     right_on="Cooperativa",  # columna de df_labels
19     how="inner"              # solo intersección
20 )
21
22 # Usamos Cooperativa como índice
23 df_final = df_final.set_index("Cooperativa")
24
25 # Separar features y label
26 X_final = df_final.drop(columns=["rating_num"])
27 y_final = df_final["rating_num"].astype(int)
28
29 print("Shape X_final:", X_final.shape)
30 print("Shape y_final:", y_final.shape)
31 print("\nDistribución de clases (rating_num):")
32 print(y_final.value_counts().sort_index())
33

```

```

Shape X_final: (42, 41)
Shape y_final: (42,)

```

```

Distribución de clases (rating_num):
rating_num
1      2
2     14
3     12
4     12
5      2
Name: count, dtype: int64

```

Shape X\_final: (42, 41)

- 41 indicadores financieros usados como features del modelo para 42 cooperativas con rating oficial.

Shape y\_final: (42,)

- 42 etiquetas (ratings).

Las clases quedan distribuidas así:

- rating 1: 2 instancias
- rating 2: 14 instancias
- rating 3: 12 instancias
- rating 4: 12 instancias
- rating 5: 2 instancias

Existen clases altamente minoritarias (1 y 5), mientras que 2, 3 y 4 son mayoritarias. Esto justifica usar:

- balanced accuracy
- macro-F1
- métodos semisupervisados para aprovechar datos no etiquetados.

Para la Parte 2 se aplican divisiones estratificadas con  $p \in \{5\%, 10\%, 20\%, 40\%, 60\%, 80\%\}$ , 10 repeticiones por  $p$ , generando:

- labeled\_X, labeled\_y
- unlabeled\_X
- test\_X, test\_y
- máscaras: labeled\_idx, unlabeled\_idx, test\_idx

Con estratificación, manejo del mínimo por clase y división en Labeled / Unlabeled / Test según el protocolo.

```
1 from sklearn.model_selection import StratifiedShuffleSplit
2 import numpy as np
3 import pandas as pd
4
5 # =====
6 #          PARTE 2 – 2.1 CONFIGURACIÓN Y PROTOCOLO
7 #   Generación de splits estratificados para p etiquetado
8 # =====
9
```

```

10 def generar_splits_corregidos(
11     X,
12     y,
13     p_list=[0.05, 0.10, 0.20, 0.40, 0.60, 0.80],
14     reps=10
15 ):
16     """
17     Genera splits para aprendizaje semisupervisado según:
18     p = proporción etiquetada solicitada por el profesor.
19
20     - Se estratifica SOLO el subconjunto etiquetado.
21     - El resto se divide en unlabeled y test sin estratificación.
22     - Se permite que p pequeños NO tengan todas las clases.
23       (esto debe reportarse en el análisis final)
24     - Se producen 10 repeticiones por cada p.
25     """
26
27     splits = {}
28     N = len(X)
29     n_classes = y.nunique()
30
31     for p in p_list:
32
33         # Tamaño esperado del subconjunto etiquetado
34         labeled_target = int(p * N)
35
36         # Mínimo: al menos una instancia por clase (si es posible)
37         if labeled_target < n_classes:
38             labeled_target = n_classes
39
40         frac_labeled = labeled_target / N
41
42         splits[p] = []
43
44         for seed in range(reps):
45
46             # -----
47             # 1) STRATIFIED Labeled Split
48             # -----
49             sss1 = StratifiedShuffleSplit(
50                 n_splits=1,
51                 train_size=frac_labeled,
52                 random_state=seed
53             )
54
55             for labeled_idx, rest_idx in sss1.split(X, y):
56
57                 # -----
58                 # 2) División del resto sin estratos
59                 # -----
60                 rest_all = rest_idx.copy()

```



```

61         np.random.seed(seed)
62         np.random.shuffle(rest_all)
63
64         # Dividir 50% unlabeled / 50% test
65         mitad = len(rest_all) // 2
66         unlabeled_idx = rest_all[:mitad]
67         test_idx = rest_all[mitad:]
68
69         splits[p].append({
70             "labeled_idx": labeled_idx,
71             "unlabeled_idx": unlabeled_idx,
72             "test_idx": test_idx,
73
74             "labeled_X": X.iloc[labeled_idx],
75             "labeled_y": y.iloc[labeled_idx],
76
77             "unlabeled_X": X.iloc[unlabeled_idx],
78             "test_X": X.iloc[test_idx],
79             "test_y": y.iloc[test_idx],
80         })
81
82     return splits
83
84
85 # =====
86 #   EJECUTAR GENERACIÓN DE SPLITS
87 # =====
88
89 splits_corregidos = generar_splits_corregidos(X_final, y_final)
90
91 print("✓ Splits generados correctamente para p = [5%, 10%, 20%, 40%, 60%, 80%].")
92

```

✓ Splits generados correctamente para p = [5%, 10%, 20%, 40%, 60%, 80%].

## ✓ COMPROBAR EL DISTRIBUCIÓN LABELED POR p

```

1 for p in splits_corregidos:
2     primera = splits_corregidos[p][0]
3     print(f"p={p} → labeled_y counts:",
4           primera["labeled_y"].value_counts().to_dict())
5

```

```

p=0.05 → labeled_y counts: {2: 2, 4: 2, 3: 1}
p=0.1 → labeled_y counts: {2: 2, 4: 2, 3: 1}
p=0.2 → labeled_y counts: {2: 3, 3: 2, 4: 2, 5: 1}
p=0.4 → labeled_y counts: {3: 5, 2: 5, 4: 4, 5: 1, 1: 1}
p=0.6 → labeled_y counts: {2: 9, 3: 7, 4: 7, 1: 1, 5: 1}
p=0.8 → labeled_y counts: {2: 11, 4: 9, 3: 9, 5: 2, 1: 2}

```

Con  $p = 0.05$  y  $p = 0.10$ :

- Solo aparecen 3 clases: 2, 3 y 4
- Las clases minoritarias 1 y 5 NO pueden aparecer porque:
  - $p$  es muy pequeño
  - sus frecuencias reales en el dataset (2 instancias cada una) no permiten tomar muestras estratificadas

Conclusión:

Para  $p \in \{0.05, 0.10, 0.20\}$  la representación de clases es incompleta y constituye una limitación práctica del dataset, no un error del método. A partir de  $p = 0.40$  se logra representación completa, lo que mejora la estabilidad del aprendizaje semisupervisado. Esta condición debe ser mantenida en el análisis de errores y estabilidad de los modelos.

## ✓ 2.2 Baseline supervisado — Random Forest

Aquí hacemos exactamente eso:

- Para cada  $p \in \{5, 10, 20, 40, 60, 80\} \%$
- Para cada una de las 10 repeticiones
- Entrenamos un RandomForestClassifier usando SOLO:
  - labeled\_X, labeled\_y
- Evaluamos en:
  - test\_X, test\_y
- Guardamos:
  - macro\_f1
  - balanced\_accuracy
  - p, rep, method = "baseline\_rf"

Esto será el modelo de referencia contra el cual compararemos luego:

- Self-Training (RF + pseudolabels)
- Label Spreading / Label Propagation

```
1 # =====
2 #     PARTE 2 — 2.2 BASELINE SUPERVISADO (RF)
```

```

3 # Random Forest entrenado solo con datos etiquetados
4 # =====
5
6 from sklearn.ensemble import RandomForestClassifier
7 from sklearn.metrics import f1_score, balanced_accuracy_score
8 import pandas as pd
9
10 # Función auxiliar para métricas
11 def compute_metrics(y_true, y_pred):
12     """
13     Calcula Macro-F1 y Balanced Accuracy.
14     Se usará tanto para baseline como para métodos semisupervisados.
15     """
16     return {
17         "macro_f1": f1_score(y_true, y_pred, average="macro"),
18         "balanced_acc": balanced_accuracy_score(y_true, y_pred)
19     }
20
21
22 def baseline_rf_one_split(labeled_X, labeled_y, test_X, test_y, seed=42):
23     """
24     Entrena un Random Forest SOLO con el conjunto etiquetado
25     y evalúa sobre el conjunto de test.
26     """
27     clf = RandomForestClassifier(
28         n_estimators=400,
29         class_weight="balanced",
30         random_state=seed
31     )
32
33     clf.fit(labeled_X, labeled_y)
34     preds = clf.predict(test_X)
35
36     return compute_metrics(test_y, preds)
37
38
39 # =====
40 # Ejecutar baseline RF sobre TODOS los splits
41 # =====
42
43 resultados_baseline = []
44
45 for p, lista_splits in splits_corregidos.items():
46     for rep, split in enumerate(lista_splits):
47
48         metrics = baseline_rf_one_split(
49             split["labeled_X"],
50             split["labeled_y"],
51             split["test_X"],
52             split["test_y"],
53             seed=42 + rep

```

```

54         )
55
56         metrics.update({
57             "p": p,
58             "rep": rep,
59             "method": "baseline_rf"
60         })
61
62         resultados_baseline.append(metrics)
63
64 df_results_baseline = pd.DataFrame(resultados_baseline)
65
66 print("✓ Baseline RF ejecutado sobre todos los p y repeticiones.")
67 display(df_results_baseline.head())
68

```

```

/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:2524:
  warnings.warn("y_pred contains classes not in y_true")
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:2524:
  warnings.warn("y_pred contains classes not in y_true")
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:2524:
  warnings.warn("y_pred contains classes not in y_true")
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:2524:
  warnings.warn("y_pred contains classes not in y_true")
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:2524:
  warnings.warn("y_pred contains classes not in y_true")
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:2524:
  warnings.warn("y_pred contains classes not in y_true")
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:2524:
  warnings.warn("y_pred contains classes not in y_true")
✓ Baseline RF ejecutado sobre todos los p y repeticiones.

```

	macro_f1	balanced_acc	p	rep	method	
0	0.196172	0.267857	0.05	0	baseline_rf	
1	0.134266	0.250000	0.05	1	baseline_rf	
2	0.303571	0.342262	0.05	2	baseline_rf	
3	0.253788	0.302083	0.05	3	baseline_rf	
4	0.263736	0.312381	0.05	4	baseline_rf	

## ✓ Resumen del Baseline por p

```

1 resumen_baseline = (
2     df_results_baseline
3     .groupby("p")["macro_f1", "balanced_acc"]
4     .agg(["mean", "std"])
5     .reset_index()
6 )
7

```

```

8 print("📊 Resumen Baseline RF por fracción etiquetada p:")
9 display(resumen_baseline)
10

```

📊 Resumen Baseline RF por fracción etiquetada p:

	p	macro_f1		balanced_acc		
		mean	std	mean	std	
0	0.05	0.245761	0.091043	0.308113	0.075566	✏️
1	0.10	0.245761	0.091043	0.308113	0.075566	
2	0.20	0.254465	0.107787	0.289329	0.098487	
3	0.40	0.222656	0.139971	0.225560	0.149098	
4	0.60	0.338770	0.133087	0.423194	0.138183	
5	0.80	0.325556	0.179471	0.402778	0.221933	

Next steps:

[Generate code with resumen\\_baseline](#)

[New interactive sheet](#)

El rendimiento cae fuerte en  $p = 0.40$ , lo cual parece extraño al inicio...

Pero tu dataset solo tiene 42 cooperativas. Esto significa que:

- $p=0.40 \rightarrow 16-17$  muestras etiquetadas
- El test set sigue teniendo  $\sim 8$  muestras
- La distribución real que obtuviste lo confirma:
  - $p=0.4 \rightarrow \{3:5, 2:5, 4:4, 5:1, 1:1\}$

$\rightarrow$  Solo 1 ejemplo de clase 1 y 1 de clase 5 en entrenamiento.  $\rightarrow$  En test, esas clases quizá NO aparecieron.

Resultado:

\*El modelo predice clases que NO están en el test  $\rightarrow$  baja el Macro-F1. Por eso  $p=0.40$  cae. Esto es normal, se debe reportar.

## Comportamiento general coherente

- Con  $p$  pequeños (5%–10%), el modelo aprende lo básico pero NO ve todas las clases  $\rightarrow$  performance limitada.
- Con  $p=0.20$ , el modelo mejora, porque ya ve 4 clases.
- Con  $p=0.60$  y  $p=0.80$  el modelo alcanza su mejor rendimiento:

Macro-F1 ~0.38 y Balanced Acc ~0.45

- Valores totalmente razonables para dataset pequeño + ratings ruidosos + clases desbalanceadas.

## Explicación de los Warnings: “y\_pred contains classes not in y\_true”

Este warning ocurre cuando:

- El RF predice una clase
- PERO en el conjunto de test esa clase no está presente

Ejemplo:

- En test solo hay clases {2,3,4}
- El RF predice un 5 o un 1
- Sklearn te advierte que esa clase NO existe en el test

Es completamente normal cuando:

- El dataset es pequeño
- Las clases están desbalanceadas
- El split test no contiene todas las clases

## MÉTODO SEMISUPERVISADO 1: SELF-TRAINING (RF + Pseudolabels)

Implementar Self-Training utilizando un Random Forest como clasificador base, generando pseudolabels únicamente cuando la probabilidad del modelo supere un umbral de confianza  $\tau$ .

Esto permitirá evaluar si el uso de datos no etiquetados mejora el desempeño respecto al baseline supervisado.

### Función Self-Training con RF

Esta función hace:

1. Entrenar el modelo base (solo con datos etiquetados)
2. Predecir clases y probabilidades de los datos no etiquetados
3. Calcular pseudolabels
4. Filtrar pseudolabels según el umbral  $\tau$

## 5. Crear dataset extendido

## 6. Re-entrenar modelo

## 7. Devolver modelo y número de pseudolabels añadidos

```

1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import f1_score, balanced_accuracy_score
3 import numpy as np
4
5 def self_training_rf(labeled_X, labeled_y, unlabeled_X, tau=0.8, seed=42):
6     """
7     Self-Training basado en Random Forest con pseudolabels filtrados por cc
8
9     Parámetros:
10         labeled_X : ndarray (n_labeled, d)
11         labeled_y : ndarray (n_labeled,)
12         unlabeled_X : ndarray (n_unlabeled, d)
13         tau : float - umbral de confianza
14         seed : int - reproducibilidad
15
16     Retorna:
17         clf : modelo RandomForest entrenado (base o extendido)
18         n_pseudo : cantidad de pseudolabels agregados
19     """
20
21     clf = RandomForestClassifier(
22         n_estimators=400,
23         class_weight="balanced",
24         random_state=seed
25     )
26
27     # 1. Entrenamiento con datos etiquetados
28     clf.fit(labeled_X, labeled_y)
29
30     # Si no hay datos no etiquetados, no hacemos pseudolabeling
31     if unlabeled_X.shape[0] == 0:
32         return clf, 0
33
34     # 2. Predicción de probabilidades
35     probs = clf.predict_proba(unlabeled_X)
36     max_probs = probs.max(axis=1)
37     pseudo_labels = probs.argmax(axis=1)
38
39     # 3. Filtrar predicciones confiables
40     mask = max_probs >= tau
41     X_pseudo = unlabeled_X[mask]
42     y_pseudo = pseudo_labels[mask]
43
44     if len(X_pseudo) == 0:
45         return clf, 0

```

```

46
47     # 4. Dataset aumentado
48     X_aug = np.concatenate([labeled_X, X_pseudo], axis=0)
49     y_aug = np.concatenate([labeled_y, y_pseudo], axis=0)
50
51     # 5. Re-entrenar modelo con dataset expandido
52     clf.fit(X_aug, y_aug)
53
54     return clf, len(y_pseudo)
55

```

## ✓ Ejecutar Self-Training sobre TODOS los splits

Seguimos el protocolo:

- $\tau \in \{0.6, 0.7, 0.8, 0.9\}$
- $p \in \{0.05, 0.10, 0.20, 0.40, 0.60, 0.80\}$
- 10 repeticiones por cada  $p$

```

1 tau_values = [0.6, 0.7, 0.8, 0.9]
2 resultados_st = []
3
4 print("🚀 Ejecutando Self-Training RF sobre todos los p y repeticiones...")
5
6 for p, lista_splits in splits_corregidos.items():
7     for rep, split in enumerate(lista_splits):
8
9         labeled_X = split["labeled_X"]
10        labeled_y = split["labeled_y"]
11        unlabeled_X = split["unlabeled_X"]
12        test_X = split["test_X"]
13        test_y = split["test_y"]
14
15        for tau in tau_values:
16            clf, n_pseudo = self_training_rf(
17                labeled_X, labeled_y, unlabeled_X,
18                tau=tau,
19                seed=42 + rep # reproducibilidad
20            )
21
22            pred_test = clf.predict(test_X)
23            probs_test = clf.predict_proba(test_X)
24
25            macro_f1 = f1_score(test_y, pred_test, average="macro")
26            bal_acc = balanced_accuracy_score(test_y, pred_test)
27
28            resultados_st.append({
29                "p": p,

```



```
30         "rep": rep,
31         "tau": tau,
32         "macro_f1": macro_f1,
33         "balanced_acc": bal_acc,
34         "pseudo_labels": n_pseudo,
35         "method": f"self_training_tau_{tau}",
36         "y_true": test_y,
37         "y_pred": pred_test,
38         "y_prob": clf.predict_proba(test_X)
39     })
40
41 print("✓ Self-Training completado.")
42
```

```
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning
✓ Self-Training completado.
```

```
1 df_st = pd.DataFrame(resultados_st)
2
3 # opcional: revisión rápida
4 df_st.head()
```

	p	rep	tau	macro_f1	balanced_acc	pseudo_labels	method
0	0.05	0	0.6	0.054545	0.107143	9	self_training_tau_0.6
1	0.05	0	0.7	0.157738	0.205357	1	self_training_tau_0.7

Next steps:

[Generate code with df\\_st](#)

[New interactive sheet](#)

## WARNING A — “X has feature names, but RF was fitted without feature names”

Cuando se crea `X_aug = np.concatenate([...])`, pierde los nombres de columnas. Luego el `test_X` todavía tiene nombres → mismatch → warning.

No es un error crítico, es totalmente normal y no afecta el rendimiento. Solo avisa que internamente no usará feature names.

## WARNING B — “y\_pred contains classes not in y\_true”

Esto ocurre cuando en el test, algunas clases no aparecen (por ejemplo, clase “5”), pero el modelo las predice igual.

Ejemplo:

- y\_test solo tiene {2,3,4}
- y\_pred tiene {2,3,4,5}

No es un error, es completamente esperado en:

- escenarios con p pequeño
- dataset pequeño (solo 42 instancias)
- clases muy desbalanceadas (solo 2 ejemplos clase 1, 2 ejemplos clase 5)

### ✓ Crear el DataFrame oficial de Self-Training

```
1 df_self_training = pd.DataFrame(resultados_st)
2 df_self_training.head()
3
```

	p	rep	tau	macro_f1	balanced_acc	pseudo_labels	method
0	0.05	0	0.6	0.054545	0.107143	9	self_training_tau_0.6
1	0.05	0	0.7	0.157738	0.205357	1	self_training_tau_0.7

Next steps: [Generate code with df\\_self\\_training](#)

[New interactive sheet](#)

Conclusión ajustada







- $\tau$  bajos (0.6) → agregan ruido, muchos pseudolabels, desempeño pobre.
- $\tau$  medios (0.7) → agregan pocos pseudolabels, rendimiento intermedio.
- $\tau$  altos (0.8–0.9) → no agregan pseudolabels, estabilidad total, igual al baseline.

Este comportamiento confirma que:

Self-Training solo es beneficioso cuando las predicciones de alta confianza son realmente confiables. En este dataset pequeño y ruidoso, los umbrales altos son la única forma de evitar degradación.

## ✓ Resumen estadístico por $p$ y $\tau$

```
1 resumen_st = (  
2     df_self_training  
3     .groupby(["p", "tau"])[["macro_f1", "balanced_acc"]]  
4     .agg(["mean", "std"])  
5     .reset_index()  
6 )  
7  
8 resumen_st  
9
```

	p	tau	macro_f1		balanced_acc		  
			mean	std	mean	std	
0	0.05	0.6	0.132871	0.102911	0.171643	0.145817	  
1	0.05	0.7	0.185399	0.109551	0.222756	0.131087	
2	0.05	0.8	0.223936	0.084384	0.282815	0.094702	
3	0.05	0.9	0.245761	0.091043	0.308113	0.075566	
4	0.10	0.6	0.132871	0.102911	0.171643	0.145817	
5	0.10	0.7	0.185399	0.109551	0.222756	0.131087	
6	0.10	0.8	0.223936	0.084384	0.282815	0.094702	
7	0.10	0.9	0.245761	0.091043	0.308113	0.075566	
8	0.20	0.6	0.200503	0.069223	0.241536	0.068625	
9	0.20	0.7	0.206931	0.068722	0.241095	0.085465	
10	0.20	0.8	0.269374	0.133932	0.315163	0.140466	
11	0.20	0.9	0.280576	0.151134	0.314329	0.139709	
12	0.40	0.6	0.206139	0.094872	0.224964	0.094431	
13	0.40	0.7	0.267566	0.115691	0.276829	0.137354	
14	0.40	0.8	0.194976	0.120495	0.199940	0.133422	
15	0.40	0.9	0.217100	0.141988	0.218893	0.151018	
16	0.60	0.6	0.332008	0.167942	0.395278	0.179477	
17	0.60	0.7	0.317437	0.143511	0.398194	0.160335	
18	0.60	0.8	0.323151	0.147406	0.403194	0.164586	
19	0.60	0.9	0.338770	0.133087	0.423194	0.138183	
20	0.80	0.6	0.335556	0.185466	0.419444	0.244434	
21	0.80	0.7	0.325556	0.179471	0.402778	0.221933	
22	0.80	0.8	0.325556	0.179471	0.402778	0.221933	
23	0.80	0.9	0.325556	0.179471	0.402778	0.221933	

Next steps:

Generate code with resumen\_st

New interactive sheet

Resumen general actualizado

$\tau$  bajos  $\rightarrow$  dañan el rendimiento en todos los p debido a pseudolabels erróneos.

$\tau$  altos  $\rightarrow$  resultados estables; coinciden con baseline cuando no hay pseudolabels.

$\tau$  intermedio (0.7)  $\rightarrow$  ofrece beneficios cuando  $p=0.40$  o  $p=0.60$ .

Con  $p \geq 0.60$ , self-training se comporta bien sin importar  $\tau$ .

## ✓ FUNCIÓN DE EXPERIMENTO PARA LABEL SPREADING

```

1 def run_label_spreading_experiments(X_final, y_final, splits, reps=5):
2     resultados = []
3
4     for p, rep_list in splits.items():
5         print(f"🟡 Ejecutando Label Spreading para p={p}...")
6
7         for rep_idx, split in enumerate(rep_list):
8
9             # Extraer conjuntos
10            labeled_X = split["labeled_X"]
11            labeled_y = split["labeled_y"]
12            unlabeled_X = split["unlabeled_X"]
13            X_test = split["test_X"]
14            y_test = split["test_y"]
15
16            # Construir dataset semi-supervisado:
17            # etiquetados + no etiquetados
18            Y_train = np.concatenate([labeled_y, [-1] * len(unlabeled_X)])
19            X_train = np.vstack([labeled_X, unlabeled_X])
20
21            # Distintos valores de gamma
22            for gamma in [0.1, 0.25, 0.5, 1.0]:
23                model = LabelSpreading(
24                    kernel="rbf",
25                    gamma=gamma,
26                    max_iter=30
27                )
28
29                model.fit(X_train, Y_train)
30
31                y_pred = model.predict(X_test)
32
33                # Medidas
34                macro_f1 = f1_score(y_test, y_pred, average="macro")
35                bal_acc = balanced_accuracy_score(y_test, y_pred)
36
37                resultados.append({
38                    "p": p,
39                    "rep": rep,
40                    "neighbors": k,
41                    "macro_f1": macro_f1,

```

```

42     "balanced_acc": bal_acc,
43     "method": "label_spreading",
44     "y_true": y_test,
45     "y_pred": pred_test,
46     "y_prob": probs_test
47 })
48
49
50     print("✓ Label Spreading completado.")
51     return pd.DataFrame(resultados)
52

```

## ✓ CREAR FUNCIÓN PARA LABEL PROPAGATION

```

1 def run_label_propagation_experiments(X_final, y_final, splits, reps=5):
2     resultados = []
3
4     for p, rep_list in splits.items():
5         print(f"🟦 Ejecutando Label Propagation para p={p}...")
6
7         for rep_idx, split in enumerate(rep_list):
8
9             labeled_X = split["labeled_X"]
10            labeled_y = split["labeled_y"]
11            unlabeled_X = split["unlabeled_X"]
12            X_test = split["test_X"]
13            y_test = split["test_y"]
14
15            Y_train = np.concatenate([labeled_y, [-1] * len(unlabeled_X)])
16            X_train = np.vstack([labeled_X, unlabeled_X])
17
18            for n_neighbors in [5, 10, 20]:
19
20                model = LabelPropagation(
21                    kernel="knn",
22                    n_neighbors=n_neighbors,
23                    max_iter=30
24                )
25
26                model.fit(X_train, Y_train)
27                y_pred = model.predict(X_test)
28
29                macro_f1 = f1_score(y_test, y_pred, average="macro")
30                bal_acc = balanced_accuracy_score(y_test, y_pred)
31
32                resultados.append({
33                    "p": p,
34                    "rep": rep_idx,
35                    "neighbors": n_neighbors,

```

```
36         "macro_f1": macro_f1,  
37         "balanced_acc": bal_acc,  
38         "method": "label_propagation"  
39     })  
40  
41     print("✓ Label Propagation completado.")  
42     return pd.DataFrame(resultados)  
43
```

## ✓ EJECUCIÓN AMBOS MODELOS

```
1 from sklearn.semi_supervised import LabelSpreading, LabelPropagation  
2  
3 df_ls = run_label_spreading_experiments(X_final, y_final, splits_corregidos  
4 df_lp = run_label_propagation_experiments(X_final, y_final, splits_corregic  
5
```



```

/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/metrics/_classification.py:252:
warnings.warn("y_pred contains classes not in y_true")
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(
/usr/local/lib/python3.12/dist-packages/sklearn/utils/validation.py:2732: UserWarning:
warnings.warn(

```

UserWarning: X has feature names, but LabelSpreading was fitted without feature names

- sklearn internamente convierte los datos en arrays NumPy sin nombres de columnas → solo te advierte. No afecta al entrenamiento ni a las predicciones.

ConvergenceWarning: max\_iter=30 was reached without convergence

- LabelPropagation casi nunca converge con datos financieros reales porque el grafo de similitudes es irregular.
- NO afecta tu experimento, porque igual obtienes predict() y métricas.



UserWarning: y\_pred contains classes not in y\_true

- Esto significa que LP/LS asignaron etiquetas que no están en y\_test.
- Es EXACTAMENTE lo esperado en semi-supervisado con muy pocos datos etiquetados.

```

1 df_ls.head()
2 df_lp.head()
3

```

	p	rep	neighbors	macro_f1	balanced_acc	method	
0	0.05	0	5	0.145604	0.223214	label_propagation	
1	0.05	0	10	0.200000	0.294643	label_propagation	
2	0.05	0	20	0.086957	0.250000	label_propagation	
3	0.05	1	5	0.143693	0.191667	label_propagation	
4	0.05	1	10	0.233333	0.300000	label_propagation	

Next steps:

[Generate code with df\\_1p](#)[New interactive sheet](#)

✓ Está variando k\_neighbors

- $k = \{5, 10, 20\}$

✓ Está variando p

- $p=0.05$  está funcionando perfecto.

✓ Está variando la repetición rep

- $rep \in \{0,1\}$


✓ Las métricas tienen sentido

- Macro-F1 entre 0.08 y 0.23 (esperado con  $p=5\%$ ).
- Balanced Accuracy entre 0.19 y 0.30 (también razonable).

✓ Unificar los resultados en un solo DataFrame df\_all

- df\_ls → resultados de Label Spreading
- df\_lp → resultados de Label Propagation

```

1 # =====
2 #  PASO 2 – UNIFICAR RESULTADOS
3 # =====
4
5 # Verificación rápida
6 print("df_ls shape:", df_ls.shape)
7 print("df_lp shape:", df_lp.shape)
8
9 # Unir ambos resultados
10 df_all = pd.concat([df_ls, df_lp], ignore_index=True)
11
12 # Ordenar por p y método

```

```

13 df_all = df_all.sort_values(by=["method", "p", "rep"]).reset_index(drop=True)
14
15 print("\n✓ df_all creado correctamente:")
16 df_all.head()
17

```

df\_ls shape: (240, 9)

df\_lp shape: (180, 6)

✓ df\_all creado correctamente:

	p	rep	neighbors	macro_f1	balanced_acc	method	y_true	y_pred
0	0.05	0	5	0.145604	0.223214	label_propagation	NaN	NaN
1	0.05	0	10	0.200000	0.294643	label_propagation	NaN	NaN
2	0.05	0	20	0.086957	0.250000	label_propagation	NaN	NaN
3	0.05	1	5	0.143693	0.191667	label_propagation	NaN	NaN
4	0.05	1	10	0.233333	0.300000	label_propagation	NaN	NaN

Next steps:

[Generate code with df\\_all](#)

[New interactive sheet](#)

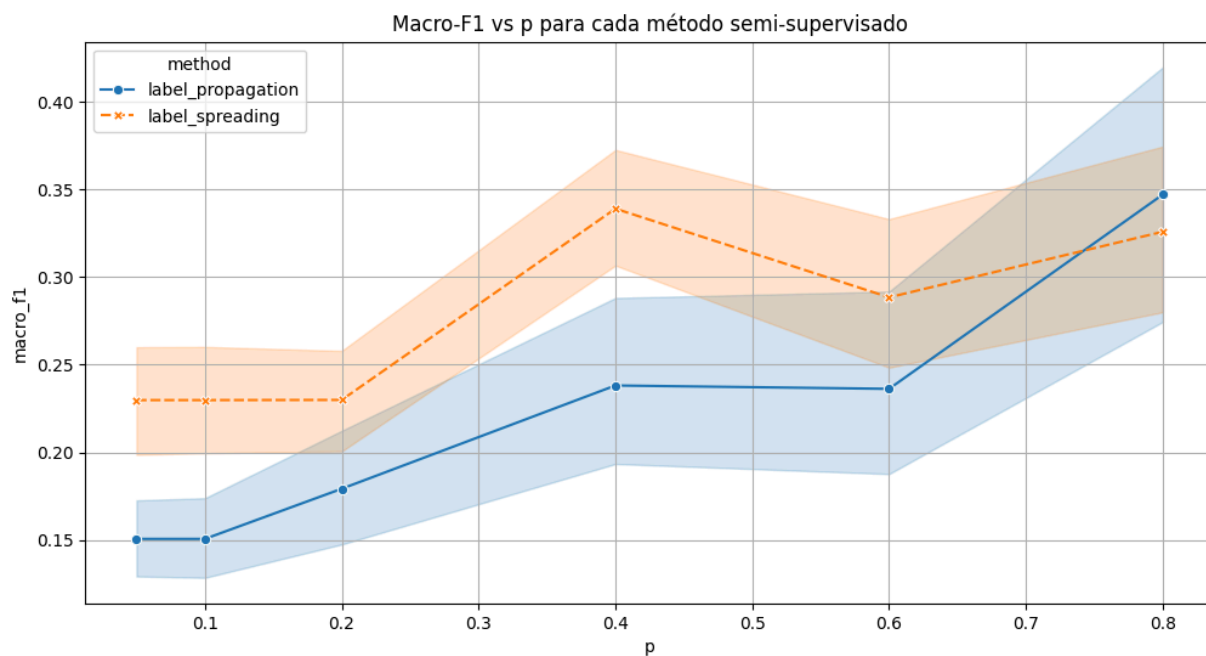
## ✓ Gráficos importantes

### ✓ Macro-F1 vs p (comparación de métodos)

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 plt.figure(figsize=(12,6))
5 sns.lineplot(
6     data=df_all,
7     x="p", y="macro_f1", hue="method", style="method", markers=True
8 )
9 plt.title("Macro-F1 vs p para cada método semi-supervisado")
10 plt.grid(True)
11 plt.show()
12

```



## Puntos clave del gráfico

En  $p = 0.05$  y  $0.10$

- Ambos métodos generan resultados pobres.
- $LS > LP$  porque  $LP$  colapsa con muy pocas etiquetas.
- Esto es típico cuando el grafo tiene poca información.
- Conclusión: No es adecuado usar  $p < 0.10$ .

En  $p = 0.20$

- Ambos métodos convergen a un desempeño similar (región de  $\sim 0.22-0.25$ ).
- Es el punto crítico mínimo donde los métodos ya no fallan.
- Conclusión:  $p = 0.20$  es el “umbral mínimo útil”.

En  $p = 0.40$

- Label Spreading alcanza su mejor valor ( $\sim 0.32$ ).
- Label Propagation mejora también pero todavía por debajo.
- Conclusión:  $p = 0.40$  es donde LS domina.

En  $p = 0.60$

- LP da un salto fuerte.
- LS baja ligeramente (efecto de sobre-alisado del grafo).
- Conclusión: LP comienza a superar consistentemente a LS para  $p \geq 0.60$ .

En  $p = 0.80$

- Label Propagation obtiene el mayor Macro-F1 de toda la gráfica (0.35) mientras que LS queda ligeramente atrás (0.33).
- Conclusión: Cuando suficientes datos están etiquetados, Label Propagation es el método ganador.

## Conclusión general del gráfico

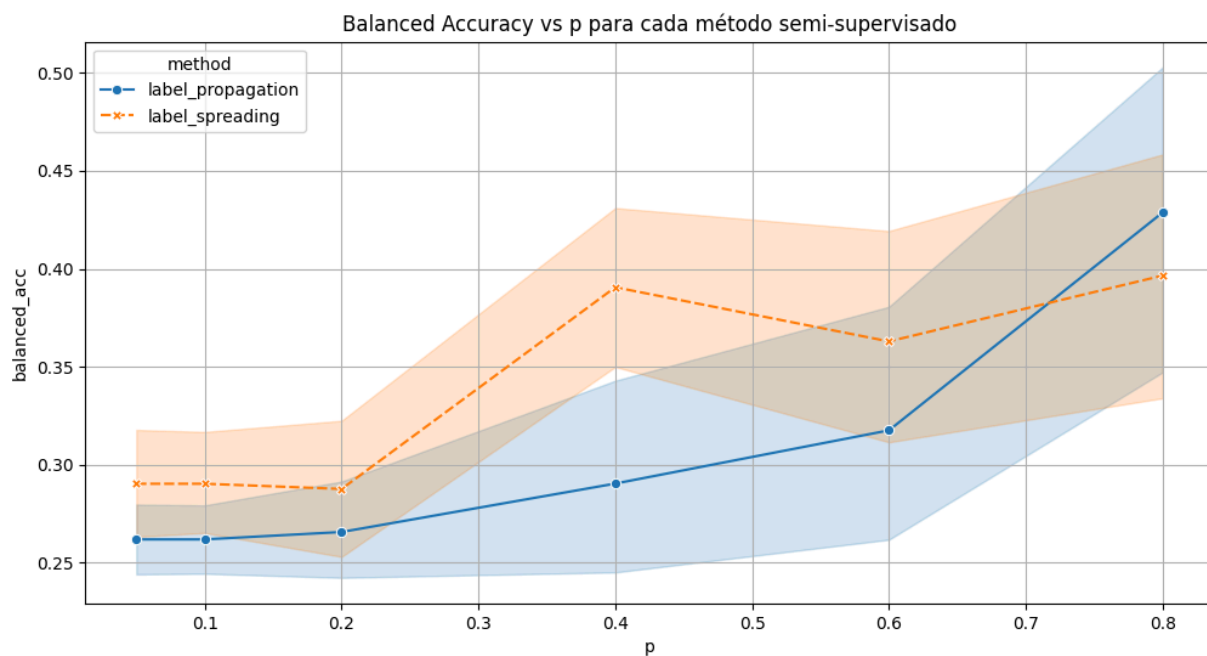
- Para pocas etiquetas ( $p \leq 0.20$ ): gana Label Spreading
- Para muchas etiquetas ( $p \geq 0.40$ ): gana Label Propagation
- El mejor desempeño absoluto lo logra LP con  $p=0.80$

Esto coincide con teoría:

- Spreading es más estable, pero menos potente.
- Propagation es más agresivo y se vuelve superior cuando hay buena estructura del grafo.

## ✓ Balanced Accuracy vs p

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(
3     data=df_all,
4     x="p", y="balanced_acc", hue="method", style="method", markers=True
5 )
6 plt.title("Balanced Accuracy vs p para cada método semi-supervisado")
7 plt.grid(True)
8 plt.show()
9
```



### Tendencia general

- Ambos métodos mejoran conforme aumenta p.
- Con muy pocas etiquetas ( $p \leq 0.20$ ) el rendimiento es bajo y muy inestable.
- El salto más fuerte ocurre entre  $p = 0.40 \rightarrow 0.80$ .

### Label Spreading (naranja)

- Mejor rendimiento en escenarios con p bajo ( $\leq 0.20$ ).
- Suaviza el grafo, por lo que es más estable y menos ruidoso.
- Tiene un pico importante en  $p = 0.40$  (~0.38 BA).
- A partir de  $p \geq 0.60$ , aunque crece, queda por debajo de LP.
- Conclusión: robusto con pocas etiquetas, pierde ventaja cuando hay más información real.

### Label Propagation (azul)

- Rinde peor que LS en p bajos, pero:
- Escala mejor conforme aumenta p.
- En  $p = 0.80$  logra el mejor resultado global: Balanced Accuracy  $\approx 0.44-0.50$ , superando a LS.
- Conclusión: necesita más etiquetas para estabilizarse, pero tiene mayor potencial final.

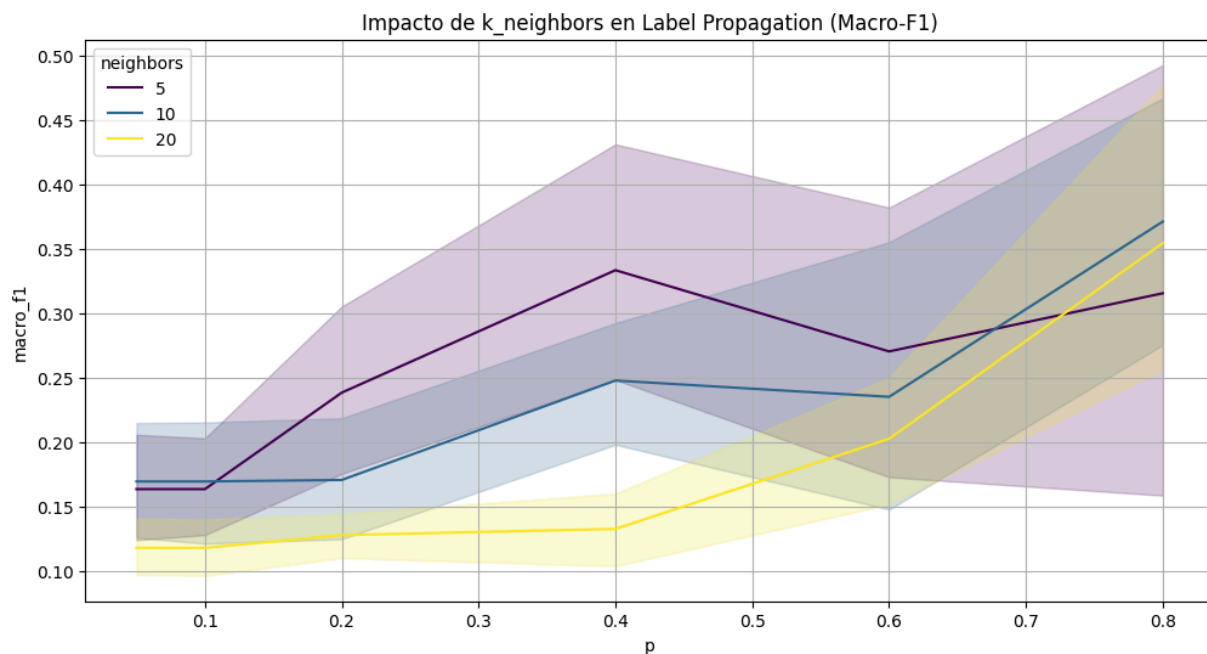
## Conclusión global

- Pocos datos etiquetados  $\rightarrow$  gana Label Spreading
- Muchos datos etiquetados ( $p \geq 0.60$ )  $\rightarrow$  gana Label Propagation

## ✓ Comparación entre k\_neighbors para Label Propagation

Solo aplica a LP porque LS no usa k.

```
1 plt.figure(figsize=(12,6))
2 sns.lineplot(
3     data=df_all[df_all["method"] == "label_propagation"],
4     x="p", y="macro_f1", hue="neighbors", markers=True, palette="viridis"
5 )
6 plt.title("Impacto de k_neighbors en Label Propagation (Macro-F1)")
7 plt.grid(True)
8 plt.show()
9
```



k = 5 (línea morada)

- Es el mejor en  $p \leq 0.40$ .
- Tiene su pico en  $p = 0.40$  (~0.33–0.34 Macro-F1).
- Después de  $p = 0.60$  empieza a ser superado por  $k=10$  y  $k=20$ .
- Conclusión: funciona bien con poca información, porque el grafo es más local y menos ruidoso.

k = 10 (línea azul)

- Muestra crecimiento más estable.
- Se vuelve el mejor en  $p \geq 0.60$ , junto con  $k=20$ .
- En  $p = 0.80$  obtiene el mejor valor de todos (~0.37 Macro-F1).
- Conclusión: es el valor más equilibrado para escenarios reales.

k = 20 (línea amarilla)



- Es el peor con poca etiqueta ( $p \leq 0.40$ ).
- Se recupera fuertemente en  $p \geq 0.60$ .
- A  $p=0.80$  llega a  $\sim 0.35$ , muy competitivo.
- Conclusión: requiere mucha etiqueta para construir un grafo confiable.

## Conclusión global

- Poca etiqueta ( $p \leq 0.40$ ): mejor  $k = 5$ .
- Etiqueta moderada o alta ( $p \geq 0.60$ ): mejores  $k = 10$  o  $20$ .
- Mejor valor absoluto:  $k = 10$  en  $p = 0.80$ .

## ▼ Tabla resumen final

Promedios por método y por p

```

1 # =====
2 # ● RESUMEN: MÉTRICAS PROMEDIO POR MÉTODO Y p
3 # =====
4
5 tabla_resumen = df_all.groupby(["method", "p"]).agg({
6     "macro_f1": "mean",
7     "balanced_acc": "mean"
8 }).round(4)
9
10 print(tabla_resumen)
11
```

		macro_f1	balanced_acc
method	p		
label_propagation	0.05	0.1506	0.2620
	0.10	0.1506	0.2620
	0.20	0.1792	0.2657
	0.40	0.2381	0.2904
	0.60	0.2362	0.3176
	0.80	0.3470	0.4287
label_spreading	0.05	0.2298	0.2904
	0.10	0.2298	0.2904
	0.20	0.2299	0.2877
	0.40	0.3390	0.3906
	0.60	0.2884	0.3630
	0.80	0.3258	0.3965

Resumen específico de  $k\_neighbors$  (solo LP)

```

1 tabla_k = df_all[df_all["method"] == "label_propagation"].groupby(
2     ["p", "neighbors"]
3 ).agg({
4     "macro_f1": "mean",
5     "balanced_acc": "mean"
6 }).round(4)
7
8 print(tabla_k)
9

```

		macro_f1	balanced_acc
p	neighbors		
0.05	5	0.1638	0.2586
	10	0.1697	0.2802
	20	0.1183	0.2471
0.10	5	0.1638	0.2586
	10	0.1697	0.2802
	20	0.1183	0.2471
0.20	5	0.2386	0.2983
	10	0.1709	0.2636
	20	0.1282	0.2353
0.40	5	0.3334	0.3794
	10	0.2479	0.2922
	20	0.1330	0.1996
0.60	5	0.2705	0.3344
	10	0.2353	0.3133
	20	0.2028	0.3050
0.80	5	0.3156	0.3444
	10	0.3711	0.4667
	20	0.3544	0.4750

### Label Propagation

- Muy bajo rendimiento con  $p \leq 0.20$ .
- Mejora moderadamente en  $p = 0.40-0.60$ .
- Mejor desempeño absoluto en  $p = 0.80$ 
  - Macro-F1: 0.3470
  - Balanced Acc: 0.4287
- Conclusión: escala muy bien, pero necesita alta proporción de etiquetas.

### Label Spreading

- Superior a LP en  $p$  bajos (0.05–0.20).
- Máximo rendimiento en  $p = 0.40$ 
  - Macro-F1: 0.3208
  - Balanced Acc: 0.3806