```python
# -*- coding: utf-8 -*-
### SCRIPT 15 - PYTHON
# PIPELINE DESARROLLO DE MODELOS DE PREDICCION, CLUSTERS, kNN Y CENTROIDES
#
# ================================================================
#
# ==========================
# Celda A+C — PRE (TRAIN-only) + Centroides + POST FULL (TRAIN & ALL)
# ==========================
# BASE: /content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2
# Entradas:
#   - <BASE>/inputs/df_test_train_v8.xlsx
# Salidas (PROD):
#   - <BASE>/modelos/preprocess_pipeline.pkl
#   - <BASE>/modelos/centroides_pre.json
# Salidas (VALID):
#   - <BASE>/validacion/01_PRE_CENTROIDES/outputs/centroides_pre_train_matrix.csv
#   - <BASE>/validacion/01_PRE_CENTROIDES/outputs/centroides_post_train_full.csv
#   - <BASE>/validacion/01_PRE_CENTROIDES/outputs/centroides_post_all_full.csv
#   - <BASE>/validacion/01_PRE_CENTROIDES/README_STEP.txt

import os, json, joblib
import numpy as np
import pandas as pd
from datetime import datetime

from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.impute import SimpleImputer

# Montaje Drive
try:
    from google.colab import drive  # type: ignore
    drive.mount('/content/drive')
except Exception:
    pass

# === Rutas base (TODO bajo _v2) ===
BASE_DIR   = "/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2"
ORIGEN_DIR = BASE_DIR
DIR_MODELOS = f"{BASE_DIR}/modelos"
DIR_VALID  = f"{BASE_DIR}/validacion"
DIR_COMP   = f"{BASE_DIR}/comparativa"  # opcional (no escribimos por defecto)

STEP_TAG = "01_PRE_CENTROIDES"
STEP_DIR = os.path.join(DIR_VALID, STEP_TAG)
OUT_DIR  = os.path.join(STEP_DIR, "outputs")
os.makedirs(DIR_MODELOS, exist_ok=True)
os.makedirs(OUT_DIR,     exist_ok=True)

random_state = 42

print(f"🏁 BASE_DIR: {BASE_DIR}")
print(f"📁 DIR_MODELOS: {DIR_MODELOS}")
print(f"🖋 DIR_VALID:   {DIR_VALID}")

# ==========
# 1) Datos
# ==========
df_path_candidates = [
    f"{ORIGEN_DIR}/inputs/df_test_train_v8.xlsx",
    f"{ORIGEN_DIR}/df_test_train_v8.xlsx",
]
df_path = next((p for p in df_path_candidates if os.path.exists(p)), None)
assert df_path is not None, f"No se encontró df_test_train_v8.xlsx en: {df_path_candidates}"
df = pd.read_excel(df_path)
print(f"✓ df_test_train_v8.xlsx cargado: {df.shape}")

# Checks mínimos
if 'Cluster_6' not in df.columns:
```

```python
 72            raise ValueError("✗ Falta 'Cluster_6'.")
 73        if 'Dataset' not in df.columns:
 74            raise ValueError("✗ Falta 'Dataset' (TRAIN/TEST).")
 75
 76        df['Cluster_6'] = df['Cluster_6'].astype(str)
 77        df['Dataset']   = df['Dataset'].astype(str).str.upper()
 78
 79        # ==========
 80        # 2) Variables PRE
 81        # ==========
 82        pre_num = [
 83            "N_lotes","C_precio_p","N_CPV","N_clasi_empresa","Presupuesto_licitacion_lote_c",
 84            "Plazo_m_c","C_juicio_valor_p_c","Intervalo_lici_d_c",
 85        ]
 86        pre_cat = [
 87            "Tipo_de_contrato_c","Tipo_de_Administracion_c","Tipo_de_procedimiento_c",
             "Tramitacion_c",
 88            # Provincia2 se añade si falta
 89            "Tipo_ganador_lote_c","Mes_lici",
 90        ]
 91
 92        # Derivar Provincia2 (dos primeros dígitos)
 93        def _to_cp_str(x):
 94            try:
 95                return str(int(x)).zfill(5)
 96            except Exception:
 97                s = str(x); return s if s and s.lower() != 'nan' else None
 98
 99        if "Provincia2" not in df.columns:
100            if "Codigo_Postal_c" in df.columns:
101                prov = df["Codigo_Postal_c"].apply(_to_cp_str).str[:2]
102                df["Provincia2"] = prov.where(prov.notna(), np.nan)
103            else:
104                df["Provincia2"] = np.nan
105        if "Provincia2" not in pre_cat:
106            pre_cat = pre_cat[:4] + ["Provincia2"] + pre_cat[4:]
107
108        # Completar columnas ausentes
109        for c in pre_num:
110            if c not in df.columns: df[c] = np.nan
111        for c in pre_cat:
112            if c not in df.columns: df[c] = np.nan
113
114        # ==========
115        # 3) PREPROCESSOR (TRAIN-only)
116        # ==========
117        df_train = df[df["Dataset"] == "TRAIN"].copy()
118        df_all   = df.copy()
119
120        num_pipe = Pipeline([
121            ("imputer", SimpleImputer(strategy="median")),
122            ("scaler",  StandardScaler())
123        ])
124        cat_pipe = Pipeline([
125            ("imputer", SimpleImputer(strategy="most_frequent")),
126            ("ohe",     OneHotEncoder(handle_unknown="ignore", sparse_output=False))
127        ])
128
129        preprocessor = ColumnTransformer(
130            transformers=[("num", num_pipe, pre_num), ("cat", cat_pipe, pre_cat)],
131            remainder='drop',
132            verbose_feature_names_out=True
133        )
134
135        X_train_pre = df_train[pre_num + pre_cat].copy()
136        y_train     = df_train["Cluster_6"].astype(str).copy()
137
138        preprocessor.fit(X_train_pre, y_train)
139        feature_names = preprocessor.get_feature_names_out().tolist()
140        print(f"• Features transformadas: {len(feature_names)}")
141
142        # ==========
```

```python
143    # 4) Centroides PRE (espacio transformado, TRAIN-only)
144    # ==========
145    Z_train = preprocessor.transform(X_train_pre)
146    classes = sorted(y_train.unique(), key=lambda x: int(x))
147
148    centroids_pre, p95_intra = {}, {}
149    for c in classes:
150        Zc = Z_train[(y_train == c).values]
151        mu = Zc.mean(axis=0)
152        centroids_pre[c] = mu.tolist()
153        d = np.sqrt(((Zc - mu) ** 2).sum(axis=1))
154        p95_intra[c] = float(np.percentile(d, 95))
155
156    # Guardar artefactos (PROD)
157    joblib.dump(preprocessor, f"{DIR_MODELOS}/preprocess_pipeline.pkl")
158    with open(f"{DIR_MODELOS}/centroides_pre.json", "w", encoding="utf-8") as f:
159        json.dump({
160            "classes": classes,
161            "feature_names": feature_names,
162            "centroids_pre": centroids_pre,
163            "p95_intra": p95_intra,
164            "random_state": random_state,
165            "built_at": datetime.now().isoformat()
166        }, f, ensure_ascii=False, indent=2)
167    print("📄 Guardado (PROD): preprocess_pipeline.pkl, centroides_pre.json")
168
169    # Copias de validación
170    cent_pre_mat = pd.DataFrame.from_dict(centroids_pre, orient='index', columns=
       feature_names)
171    cent_pre_mat.index.name = "Cluster_6"
172    cent_pre_mat.sort_index(key=lambda x: x.map(int), inplace=True)
173    cent_pre_mat.to_csv(f"{OUT_DIR}/centroides_pre_train_matrix.csv")
174
175    def _top12_share(series: pd.Series):
176        vc = series.value_counts(dropna=True)
177        if len(vc) == 0: return (np.nan, np.nan, np.nan, np.nan)
178        top1, s1 = vc.index[0], float(vc.iloc[0]) / float(vc.sum())
179        if len(vc) > 1:
180            top2, s2 = vc.index[1], float(vc.iloc[1]) / float(vc.sum())
181        else:
182            top2, s2 = (np.nan, np.nan)
183        return (top1, s1, top2, s2)
184
185    def _post_full(df_base: pd.DataFrame, idx_vals):
186        pre_num_present = [c for c in pre_num if c in df_base.columns]
187        post_num = [c for c in ["N_ofertantes","Ofertas_admitidas","Desierta",
           "Importe_adjudicacion_lote"] if c in df_base.columns]
188        out = pd.DataFrame(index=idx_vals)
189        if pre_num_present:
190            g = df_base.groupby("Cluster_6")
191            pre_mean = g[pre_num_present].mean(numeric_only=True).add_prefix("PRE__").
               add_suffix("__mean")
192            pre_std  = g[pre_num_present].std(numeric_only=True).add_prefix("PRE__").
               add_suffix("__std")
193            out = out.join(pre_mean, how="left").join(pre_std, how="left")
194        if post_num:
195            g = df_base.groupby("Cluster_6")
196            post_mean = g[post_num].mean(numeric_only=True).add_prefix("POST__").
               add_suffix("__mean")
197            post_std  = g[post_num].std(numeric_only=True).add_prefix("POST__").add_suffix
               ("__std")
198            out = out.join(post_mean, how="left").join(post_std,  how="left")
199        rows = []
200        for c in idx_vals:
201            sub = df_base[df_base["Cluster_6"].astype(str) == c]
202            row = {"Cluster_6": c}
203            for col in pre_cat:
204                if col not in sub.columns:
205                    row[f"PRECAT__{col}__top1"]=np.nan; row[f"PRECAT__{col}__top1_share"]=
                       np.nan
206                    row[f"PRECAT__{col}__top2"]=np.nan; row[f"PRECAT__{col}__top2_share"]=
                       np.nan
```

```python
                    else:
                        top1, s1, top2, s2 = _top12_share(sub[col])
                        row[f"PRECAT__{col}__top1"]=top1; row[f"PRECAT__{col}__top1_share"]=s1
                        row[f"PRECAT__{col}__top2"]=top2; row[f"PRECAT__{col}__top2_share"]=s2
                rows.append(row)
        out2 = pd.DataFrame(rows).set_index("Cluster_6")
        final = out.join(out2, how="left")
        final.index.name = "Cluster_6"
        final.sort_index(key=lambda x: x.map(int), inplace=True)
        return final

    post_train = _post_full(df_train.copy(), classes)
    post_all   = _post_full(df_all.copy(),   classes)
    post_train.to_csv(f"{OUT_DIR}/centroides_post_train_full.csv")
    post_all.to_csv(f"{OUT_DIR}/centroides_post_all_full.csv")

    with open(os.path.join(STEP_DIR, "README_STEP.txt"), "w", encoding="utf-8") as f:
        f.write(f"""STEP: {STEP_TAG}
BASE_DIR: {BASE_DIR}
Fecha: {datetime.now().isoformat()}
random_state: {random_state}
Artefactos PROD: modelos/preprocess_pipeline.pkl, modelos/centroides_pre.json
Salidas VALID: outputs/*.csv
""")
    print("✅ Celda A+C COMPLETADA (todo bajo _v2).")

"""Celda D — Distancias a centroides (TRAIN/TEST)"""

# ============================
# Celda D — Distancias a centroides PRE (TRAIN/TEST)
# ============================
# Entradas:
#   - <BASE>/inputs/df_test_train_v8.xlsx
#   - <BASE>/modelos/preprocess_pipeline.pkl
#   - <BASE>/modelos/centroides_pre.json
# Salidas (VALID):
#   - <BASE>/validacion/03_DISTANCIAS_PRE/outputs/df_train_with_dist.csv
#   - <BASE>/validacion/03_DISTANCIAS_PRE/outputs/df_test_with_dist.csv
#   - <BASE>/validacion/03_DISTANCIAS_PRE/outputs/summary_distancias.json

import os, json, joblib
import numpy as np
import pandas as pd
from datetime import datetime

try:
    from google.colab import drive  # type: ignore
    drive.mount('/content/drive')
except Exception:
    pass

BASE_DIR   = "/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2"
ORIGEN_DIR = BASE_DIR
DIR_MODELOS = f"{BASE_DIR}/modelos"
DIR_VALID  = f"{BASE_DIR}/validacion"

STEP_TAG = "03_DISTANCIAS_PRE"
STEP_DIR = os.path.join(DIR_VALID, STEP_TAG)
OUT_DIR  = os.path.join(STEP_DIR, "outputs")
os.makedirs(OUT_DIR, exist_ok=True)

# Artefactos
preproc_path   = os.path.join(DIR_MODELOS, "preprocess_pipeline.pkl")
cent_json_path = os.path.join(DIR_MODELOS, "centroides_pre.json")
assert os.path.exists(preproc_path), "Falta preprocess_pipeline.pkl"
assert os.path.exists(cent_json_path), "Falta centroides_pre.json"

preprocessor = joblib.load(preproc_path)
with open(cent_json_path, "r", encoding="utf-8") as f:
    cent_meta = json.load(f)

classes        = [str(c) for c in cent_meta["classes"]]
```

```python
279    centroids_pre = {str(k): np.array(v, dtype=float) for k, v in cent_meta[
       "centroids_pre"].items()}
280
281    # Datos
282    df_path_candidates = [
283        f"{ORIGEN_DIR}/inputs/df_test_train_v8.xlsx",
284        f"{ORIGEN_DIR}/df_test_train_v8.xlsx",
285    ]
286    df_path = next((p for p in df_path_candidates if os.path.exists(p)), None)
287    assert df_path is not None, "No se encontró df_test_train_v8.xlsx"
288    df = pd.read_excel(df_path)
289    if "Cluster_6" in df.columns: df["Cluster_6"] = df["Cluster_6"].astype(str)
290    df["Dataset"] = df["Dataset"].astype(str).str.upper()
291
292    pre_num = [
293        "N_lotes","C_precio_p","N_CPV","N_clasi_empresa","Presupuesto_licitacion_lote_c",
294        "Plazo_m_c","C_juicio_valor_p_c","Intervalo_lici_d_c",
295    ]
296    pre_cat = [
297        "Tipo_de_contrato_c","Tipo_de_Administracion_c","Tipo_de_procedimiento_c",
           "Tramitacion_c",
298        "Provincia2","Tipo_ganador_lote_c","Mes_lici",
299    ]
300
301    def _to_cp_str(x):
302        try:
303            return str(int(x)).zfill(5)
304        except Exception:
305            s = str(x); return s if s and s.lower() != 'nan' else None
306
307    def _prepare_pre(df_in: pd.DataFrame) -> pd.DataFrame:
308        df = df_in.copy()
309        if "Provincia2" not in df.columns:
310            if "Codigo_Postal_c" in df.columns:
311                prov = df["Codigo_Postal_c"].apply(_to_cp_str).str[:2]
312                df["Provincia2"] = prov.where(prov.notna(), np.nan)
313            else:
314                df["Provincia2"] = np.nan
315        if "Mes_lici" not in df.columns: df["Mes_lici"] = np.nan
316        for c in pre_num:
317            if c not in df.columns: df[c] = np.nan
318        for c in pre_cat:
319            if c not in df.columns: df[c] = np.nan
320        return df[pre_num + pre_cat]
321
322    df_train = df[df["Dataset"] == "TRAIN"].copy()
323    df_test  = df[df["Dataset"] == "TEST"].copy()
324    Z_train  = preprocessor.transform(_prepare_pre(df_train))
325    Z_test   = preprocessor.transform(_prepare_pre(df_test))
326
327    def _dist_to_centroids(Z):
328        rows=[]
329        for i in range(Z.shape[0]):
330            z=Z[i]; dists={}
331            for c in classes:
332                mu = centroids_pre[c]
333                dists[f"dist_c{c}"] = float(np.sqrt(((z - mu) ** 2).sum()))
334            items = sorted(dists.items(), key=lambda x: x[1])
335            min_k, min_d = items[0]; second_d = items[1][1] if len(items)>1 else np.nan
336            rows.append({"min_cluster": min_k.replace("dist_c",""), "min_dist": min_d,
                "margin": float(second_d-min_d) if np.isfinite(second_d) else np.nan, **dists
                })
337        return pd.DataFrame(rows)
338
339    dist_train = _dist_to_centroids(Z_train)
340    dist_test  = _dist_to_centroids(Z_test)
341
342    train_out = pd.concat([df_train.reset_index(drop=True), dist_train], axis=1)
343    test_out  = pd.concat([df_test.reset_index(drop=True),  dist_test],  axis=1)
344
345    train_out_path = os.path.join(OUT_DIR, "df_train_with_dist.csv")
346    test_out_path  = os.path.join(OUT_DIR, "df_test_with_dist.csv")
```

```python
347    train_out.to_csv(train_out_path, index=False)
348    test_out.to_csv(test_out_path,  index=False)
349
350    with open(os.path.join(OUT_DIR, "summary_distancias.json"), "w", encoding="utf-8") as
       f:
351        json.dump({
352            "step": STEP_TAG,
353            "built_at": datetime.now().isoformat(),
354            "outputs": {"train_with_dist": train_out_path, "test_with_dist": test_out_path
               }
355        }, f, ensure_ascii=False, indent=2)
356
357    with open(os.path.join(STEP_DIR, "README_STEP.txt"), "w", encoding="utf-8") as f:
358        f.write(f"{STEP_TAG} generado el {datetime.now().isoformat()} (BASE_DIR={BASE_DIR
               }).\n")
359
360    print("✅ Distancias PRE generadas (todo bajo _v2).")
361
362    """Celda A — META del PRE + centroides (preprocess_pipeline)"""
363
364    # =========================================
365    # META — PREPROCESS + CENTROIDES (guardado)
366    # =========================================
367    import json, joblib, numpy as np
368    from pathlib import Path
369
370    BASE = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
371    DIR_MODELOS = BASE / "modelos"
372    DIR_MODELOS.mkdir(parents=True, exist_ok=True)
373
374    preproc_path   = DIR_MODELOS / "preprocess_pipeline.pkl"
375    centroids_path = DIR_MODELOS / "centroides_pre.json"
376    assert preproc_path.exists() and centroids_path.exists(), "Faltan
       preprocess_pipeline.pkl y/o centroides_pre.json"
377
378    pre = joblib.load(preproc_path)
379    cent = json.loads(centroids_path.read_text(encoding="utf-8"))
380
381    # Extraer columnas crudas esperadas por el ColumnTransformer
382    num_cols = []
383    cat_cols = []
384    try:
385        trfs = pre.named_transformers_
386        if "num" in trfs: num_cols = list(pre.transformers_[0][2]) if isinstance(pre.
           transformers_[0][2], list) else []
387        if "cat" in trfs: cat_cols = list(pre.transformers_[1][2]) if isinstance(pre.
           transformers_[1][2], list) else []
388    except Exception:
389        pass
390
391    pre_meta = {
392        "artifact": "preprocess_pipeline.pkl",
393        "expects_raw_features": {
394            "num": num_cols,
395            "cat": cat_cols
396        },
397        "transformers": {
398            "num": [s[0] for s in getattr(pre.named_transformers_.get("num", None),
               "steps", [])] if "num" in pre.named_transformers_ else [],
399            "cat": [s[0] for s in getattr(pre.named_transformers_.get("cat", None),
               "steps", [])] if "cat" in pre.named_transformers_ else []
400        },
401        "centroids_ref": "centroides_pre.json",
402        "centroids_overview": {
403            "classes": cent.get("classes"),
404            "p95_intra_available": bool(cent.get("p95_intra"))
405        }
406    }
407    (DIR_MODELOS / "preprocess_meta.json").write_text(json.dumps(pre_meta, indent=2,
       ensure_ascii=False), encoding="utf-8")
408    print("✅ Guardado modelos/preprocess_meta.json")
409
```

```python
"""Celda E — kNN PRE (TRAIN y ALL) — NUM-ONLY"""

# ============================
# Celda E — kNN de similares (TRAIN y ALL) — NUM-ONLY
# ============================
# Entradas:
#   - <BASE>/inputs/df_test_train_v8.xlsx
# Salidas (PROD):
#   - <BASE>/modelos/knn_pre_train.joblib
#   - <BASE>/modelos/knn_pre_all.joblib
#   - <BASE>/modelos/knn_params.json
# Salidas (VALID):
#   - <BASE>/validacion/04_KNN_PRE/outputs/knn_train_cases.parquet
#   - <BASE>/validacion/04_KNN_PRE/outputs/knn_all_cases.parquet
#   - <BASE>/validacion/04_KNN_PRE/outputs/meta_knn.json

import os, json, joblib
import numpy as np
import pandas as pd
from datetime import datetime
from sklearn.neighbors import NearestNeighbors
from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler

try:
    from google.colab import drive  # type: ignore
    drive.mount('/content/drive')
except Exception:
    pass

BASE_DIR   = "/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2"
ORIGEN_DIR = BASE_DIR
DIR_MODELOS = f"{BASE_DIR}/modelos"
DIR_VALID   = f"{BASE_DIR}/validacion"

STEP_TAG = "04_KNN_PRE"
STEP_DIR = os.path.join(DIR_VALID, STEP_TAG)
OUT_DIR  = os.path.join(STEP_DIR, "outputs")
os.makedirs(DIR_MODELOS, exist_ok=True)
os.makedirs(OUT_DIR,     exist_ok=True)

# Datos
df_path_candidates = [
    f"{ORIGEN_DIR}/inputs/df_test_train_v8.xlsx",
    f"{ORIGEN_DIR}/df_test_train_v8.xlsx",
]
df_path = next((p for p in df_path_candidates if os.path.exists(p)), None)
assert df_path is not None, "No se encontró df_test_train_v8.xlsx"
df = pd.read_excel(df_path)
df["Dataset"] = df["Dataset"].astype(str).str.upper()

df_train = df[df["Dataset"] == "TRAIN"].copy()
df_test  = df[df["Dataset"] == "TEST"].copy()
df_all   = pd.concat([df_train, df_test], axis=0, ignore_index=True)

# PRE NUM-ONLY
pre_num = [
    "N_lotes","C_precio_p","N_CPV","N_clasi_empresa","Presupuesto_licitacion_lote_c",
    "Plazo_m_c","C_juicio_valor_p_c","Intervalo_lici_d_c",
]

def _prepare_num(df_in: pd.DataFrame) -> pd.DataFrame:
    df = df_in.copy()
    for c in pre_num:
        if c not in df.columns: df[c] = np.nan
        df[c] = pd.to_numeric(df[c], errors="coerce")
    return df[pre_num]

X_train_num = _prepare_num(df_train)
X_all_num   = _prepare_num(df_all)
```

```python
482    num_only = Pipeline([
483        ("imputer", SimpleImputer(strategy="median")),
484        ("scaler",  StandardScaler())
485    ])
486    num_only.fit(X_train_num)  # TRAIN-only
487
488    Z_train = num_only.transform(X_train_num)
489    Z_all   = num_only.transform(X_all_num)
490
491    # kNN + p95
492    K_NEIGH = 25
493    nn_train = NearestNeighbors(n_neighbors=K_NEIGH, metric="euclidean").fit(Z_train)
494    nn_all   = NearestNeighbors(n_neighbors=K_NEIGH, metric="euclidean").fit(Z_all)
495
496    nn_tmp   = NearestNeighbors(n_neighbors=2, metric="euclidean").fit(Z_train)
497    dist2, _ = nn_tmp.kneighbors(Z_train, n_neighbors=2, return_distance=True)
498    nn1_dists = dist2[:, 1]
499    p95_nn1_train = float(np.percentile(nn1_dists, 95))
500    print(f"🔍 p95_nn1_train (NUM-ONLY) = {p95_nn1_train:.4f}")
501
502    # Guardar artefactos PROD
503    joblib.dump(nn_train, os.path.join(DIR_MODELOS, "knn_pre_train.joblib"))
504    joblib.dump(nn_all,   os.path.join(DIR_MODELOS, "knn_pre_all.joblib"))
505    with open(os.path.join(DIR_MODELOS, "knn_params.json"), "w", encoding="utf-8") as f:
506        json.dump({
507            "built_at": datetime.now().isoformat(),
508            "k_neighbors": K_NEIGH,
509            "p95_nn1_train": p95_nn1_train,
510            "metric": "euclidean",
511            "index_key": "knn_row",
512            "space": "pre_num",
513            "pre_num_features": pre_num
514        }, f, ensure_ascii=False, indent=2)
515    print("📄 Guardado (PROD): knn_pre_train.joblib, knn_pre_all.joblib, knn_params.json")
516
517    # Casos VALID
518    def _attach_knn_row(df_src: pd.DataFrame) -> pd.DataFrame:
519        out = df_src.reset_index(drop=True).copy()
520        out.insert(0, "knn_row", np.arange(len(out), dtype=int))
521        return out
522
523    cases_train = _attach_knn_row(df_train)
524    cases_all   = _attach_knn_row(df_all)
525
526    train_cases_path = os.path.join(OUT_DIR, "knn_train_cases.parquet")
527    all_cases_path   = os.path.join(OUT_DIR, "knn_all_cases.parquet")
528    cases_train.to_parquet(train_cases_path, index=False)
529    cases_all.to_parquet(all_cases_path, index=False)
530
531    with open(os.path.join(OUT_DIR, "meta_knn.json"), "w", encoding="utf-8") as f:
532        json.dump({
533            "artifacts": {"nn_train": "modelos/knn_pre_train.joblib", "nn_all":
                "modelos/knn_pre_all.joblib"},
534            "cases": {"train": train_cases_path, "all": all_cases_path, "format":
                "parquet", "index_key": "knn_row"},
535            "p95_nn1_train": p95_nn1_train, "space": "pre_num"
536        }, f, ensure_ascii=False, indent=2)
537
538    with open(os.path.join(STEP_DIR, "README_STEP.txt"), "w", encoding="utf-8") as f:
539        f.write(f"{STEP_TAG} (NUM-ONLY) generado el {datetime.now().isoformat()}
                (BASE_DIR={BASE_DIR}).\n")
540
541    print("☑ kNN NUM-ONLY listo (todo bajo _v2).")
542
543    """Celda B — META del kNN (incluye pipeline num-only fitted)"""
544
545    # ==========================================
546    # META — KNN (pipeline num-only + meta JSON)
547    # ==========================================
548    import json, joblib
549    from pathlib import Path
550
```

```python
551    BASE = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
552    DIR_MODELOS = BASE / "modelos"
553    assert 'num_only' in globals(), "num_only no está en memoria: ejecuta esta celda
       justo tras la Celda E (KNN)."
554
555    # Guardar pipeline num-only usada realmente
556    knn_num_pipe_path = DIR_MODELOS / "knn_num_pipeline.joblib"
557    joblib.dump(num_only, knn_num_pipe_path)
558
559    # Actualizar/crear meta del KNN (aprovecha knn_params.json creado en Celda E)
560    knn_params_path = DIR_MODELOS / "knn_params.json"
561    assert knn_params_path.exists(), "Falta modelos/knn_params.json (Célula E)."
562    params = json.loads(knn_params_path.read_text(encoding="utf-8"))
563    params["num_only_pipeline_file"] = "knn_num_pipeline.joblib"
564    params["note"] = "Pipeline num-only exacta usada para estandarizar las 8 features
       (NUM-ONLY) antes del índice."
565    knn_params_path.write_text(json.dumps(params, indent=2, ensure_ascii=False), encoding=
       "utf-8")
566
567    print("✅ Guardado modelos/knn_num_pipeline.joblib y actualizado
       modelos/knn_params.json")
568
569    """Celda F — Validación kNN (ALL vs TRAIN) con TEST"""
570
571    # ============================================
572    # Celda E+1 — Catálogo kNN (ALL) COMPLETO
573    # ============================================
574    # BASE: /content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2
575    # Entradas:
576    #   - <BASE>/inputs/df_test_train_v8.xlsx
577    #   - (opcional) <BASE>/validacion/04_KNN_PRE/outputs/knn_all_cases.parquet  (para
       verificar alineación)
578    # Salidas (PROD):
579    #   - <BASE>/modelos/knn_catalog_all.parquet        (TODAS las variables + knn_row)
580    #   - <BASE>/modelos/knn_catalog_meta.json
581    # Salidas (VALID):
582    #   - <BASE>/validacion/04_KNN_PRE/outputs/knn_catalog_all_head.csv   (muestra)
583    #   - <BASE>/validacion/04_KNN_PRE/outputs/knn_catalog_alignment.json (reporte
       verificación)
584
585    import os, json
586    import numpy as np
587    import pandas as pd
588    from datetime import datetime
589
590    # Rutas base
591    BASE_DIR   = "/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2"
592    DIR_MODELOS = f"{BASE_DIR}/modelos"
593    DIR_VALID  = f"{BASE_DIR}/validacion"
594    STEP_DIR   = os.path.join(DIR_VALID, "04_KNN_PRE")
595    OUT_DIR    = os.path.join(STEP_DIR, "outputs")
596    os.makedirs(DIR_MODELOS, exist_ok=True)
597    os.makedirs(OUT_DIR, exist_ok=True)
598
599    # 1) Cargar df y recrear ALL (TRAIN+TEST)
600    df_path_candidates = [
601        f"{BASE_DIR}/inputs/df_test_train_v8.xlsx",
602        f"{BASE_DIR}/df_test_train_v8.xlsx",
603    ]
604    df_path = next((p for p in df_path_candidates if os.path.exists(p)), None)
605    assert df_path is not None, "No se encontró df_test_train_v8.xlsx"
606    df = pd.read_excel(df_path)
607    assert "Dataset" in df.columns, "✗ Falta 'Dataset' (TRAIN/TEST)"
608    df["Dataset"] = df["Dataset"].astype(str).str.upper()
609
610    df_train = df[df["Dataset"] == "TRAIN"].copy()
611    df_test  = df[df["Dataset"] == "TEST"].copy()
612    df_all   = pd.concat([df_train, df_test], axis=0, ignore_index=True)
613
614    # 2) Asegurar Provincia2 si falta (regla del proyecto)
615    def _to_cp_str(x):
616        try: return str(int(x)).zfill(5)
```

```python
        except Exception:
            s = str(x); return s if (isinstance(s, str) and s and s.lower()!="nan") else
            np.nan

    if "Provincia2" not in df_all.columns:
        if "Codigo_Postal_c" in df_all.columns:
            df_all["Provincia2"] = df_all["Codigo_Postal_c"].apply(_to_cp_str).str[:2]
        else:
            df_all["Provincia2"] = np.nan

    # 3) Asegurar columna de ID amigable (no sobrescribimos si ya existe)
    id_candidates = ["Identificador","identificador","ID","Id","id_lote","id_contrato"]
    id_col = next((c for c in id_candidates if c in df_all.columns), None)
    if id_col is None:
        df_all["Identificador"] = [f"row_{i}" for i in range(len(df_all))]
        id_col = "Identificador"

    # 4) Insertar knn_row para alinear con Celda E (reset_index + rango continuo)
    catalog = df_all.reset_index(drop=True).copy()
    catalog.insert(0, "knn_row", np.arange(len(catalog), dtype=int))

    # 5) Guardar catálogo completo (PRODUCCIÓN)
    catalog_path = os.path.join(DIR_MODELOS, "knn_catalog_all.parquet")
    catalog.to_parquet(catalog_path, index=False)

    # 6) Muestra y metadatos (VALIDACIÓN)
    head_path = os.path.join(OUT_DIR, "knn_catalog_all_head.csv")
    catalog.head(30).to_csv(head_path, index=False)

    meta = {
        "built_at": datetime.now().isoformat(),
        "rows": int(len(catalog)),
        "cols": int(catalog.shape[1]),
        "id_col": id_col,
        "index_key": "knn_row",
        "source_df": os.path.basename(df_path),
        "note": "Catálogo ALL (TRAIN+TEST) con TODAS las variables + knn_row al frente."
    }
    with open(os.path.join(DIR_MODELOS, "knn_catalog_meta.json"), "w", encoding="utf-8")
    as f:
        json.dump(meta, f, ensure_ascii=False, indent=2)

    # 7) Verificación de alineación con knn_all_cases.parquet (si existe)
    align_report = {"checked": False, "ok": None, "details": {}}
    cases_all_path = os.path.join(OUT_DIR, "knn_all_cases.parquet")
    if os.path.exists(cases_all_path):
        cases_all = pd.read_parquet(cases_all_path)
        align_report["checked"] = True
        try:
            # mismas filas y mismo rango de knn_row
            same_len = (len(cases_all) == len(catalog))
            same_range = (cases_all["knn_row"].iloc[0] == 0 and cases_all["knn_row"].iloc
            [-1] == len(catalog)-1)
            # chequeo de identidad en 5 posiciones aleatorias
            rng = np.random.RandomState(42)
            sample_idx = rng.choice(len(catalog), size=min(5, len(catalog)), replace=False
            )
            id_match = []
            for i in sample_idx:
                r_cases = cases_all.iloc[i]
                r_cata  = catalog.iloc[i]
                id_match.append(bool(r_cases.get("knn_row", -1) == r_cata.get("knn_row", -
                2)))
            align_ok = same_len and same_range and all(id_match)
            align_report["ok"] = bool(align_ok)
            align_report["details"] = {
                "same_len": bool(same_len),
                "same_range": bool(same_range),
                "sample_id_match_rate": float(np.mean(id_match)) if len(id_match)>0 else
                None
            }
        except Exception as e:
```

```python
            align_report["ok"] = False
            align_report["details"] = {"error": str(e)}

    with open(os.path.join(OUT_DIR, "knn_catalog_alignment.json"), "w", encoding="utf-8")
    as f:
        json.dump(align_report, f, ensure_ascii=False, indent=2)

    print("☑ Catálogo kNN ALL generado.")
    print(f"→ modelos/knn_catalog_all.parquet  | filas={len(catalog)} cols={catalog.shape[
    1]}")
    print(f"→ validacion/04_KNN_PRE/outputs/knn_catalog_all_head.csv")
    if align_report["checked"]:
        print(f"✔ Verificación con knn_all_cases.parquet | ok={align_report['ok']} |
        details={align_report['details']}")
    else:
        print("ⓘ No se encontró knn_all_cases.parquet; se omitió verificación de
        alineación.")

    """Celda E+1 — Catálogo kNN (ALL) COMPLETO"""

    # ============================
    # Celda F — Validación kNN (ALL vs TRAIN) con TEST
    # ============================
    # Salidas (VALID):
    #   - <BASE>/validacion/05_KNN_VALIDACION/outputs/knn_test_validation_pairs_all.csv
    #   - <BASE>/validacion/05_KNN_VALIDACION/outputs/knn_test_validation_pairs_train.csv
    #   - <BASE>/validacion/05_KNN_VALIDACION/outputs/knn_test_hit3_by_cluster_cases.csv
    #   - <BASE>/validacion/05_KNN_VALIDACION/outputs/knn_validation_metrics_extended.json

    import os, json, joblib
    import numpy as np
    import pandas as pd
    from datetime import datetime
    from sklearn.neighbors import NearestNeighbors

    try:
        from google.colab import drive  # type: ignore
        drive.mount('/content/drive')
    except Exception:
        pass

    BASE_DIR   = "/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2"
    ORIGEN_DIR = BASE_DIR
    DIR_MODELOS = f"{BASE_DIR}/modelos"
    DIR_VALID   = f"{BASE_DIR}/validacion"

    STEP_TAG = "05_KNN_VALIDACION"
    STEP_DIR = os.path.join(DIR_VALID, STEP_TAG)
    OUT_DIR  = os.path.join(STEP_DIR, "outputs")
    os.makedirs(OUT_DIR, exist_ok=True)

    # Artefactos kNN
    nn_train_path   = os.path.join(DIR_MODELOS, "knn_pre_train.joblib")
    nn_all_path     = os.path.join(DIR_MODELOS, "knn_pre_all.joblib")
    knn_params_path = os.path.join(DIR_MODELOS, "knn_params.json")
    assert os.path.exists(nn_train_path) and os.path.exists(nn_all_path) and os.path.
    exists(knn_params_path), "Faltan artefactos kNN"

    with open(knn_params_path, "r", encoding="utf-8") as f:
        knn_params = json.load(f)
    p95_nn1_train = float(knn_params["p95_nn1_train"])
    INDEX_KEY = knn_params.get("index_key", "knn_row")

    nn_train = joblib.load(nn_train_path)
    nn_all   = joblib.load(nn_all_path)

    # Casos (de Celda E)
    cases_train_path = os.path.join(DIR_VALID, "04_KNN_PRE", "outputs",
    "knn_train_cases.parquet")
    cases_all_path   = os.path.join(DIR_VALID, "04_KNN_PRE", "outputs",
    "knn_all_cases.parquet")
    cases_train = pd.read_parquet(cases_train_path)
```

```python
748    cases_all    = pd.read_parquet(cases_all_path)
749
750    # Datos TEST
751    df_path_candidates = [
752        f"{ORIGEN_DIR}/inputs/df_test_train_v8.xlsx",
753        f"{ORIGEN_DIR}/df_test_train_v8.xlsx",
754    ]
755    df_path = next((p for p in df_path_candidates if os.path.exists(p)), None)
756    assert df_path is not None, "No se encontró df_test_train_v8.xlsx"
757    df = pd.read_excel(df_path)
758    df["Dataset"] = df["Dataset"].astype(str).str.upper()
759    df_test = df[df["Dataset"] == "TEST"].copy().reset_index(drop=True)
760
761    # Columna ID para excluir 'self' por ID
762    id_col = None
763    for cand in ["Identificador","identificador","ID","Id","id_lote","id_contrato"]:
764        if cand in df_test.columns: id_col = cand; break
765    if id_col is None:
766        df_test["Identificador"] = [f"test_row_{i}" for i in range(len(df_test))]
767        id_col = "Identificador"
768
769    def sim_percent(d, p95=p95_nn1_train):
770        return max(0.0, 1.0 - (float(d) / p95)) * 100.0 if p95 > 0 else 0.0
771
772    # Consultas (los objetos nn_* ya están ajustados sobre TRAIN/ALL en NUM-ONLY)
773    dist_all,  idx_all  = nn_all.kneighbors(n_neighbors=10,  return_distance=True, X=None)
774    dist_train,idx_train= nn_train.kneighbors(n_neighbors=10, return_distance=True, X=None
775    )
776    # Construcción de pares
777    def build_pairs_csv(df_test, idx_mat, dist_mat, cases_df, out_csv_path,
       exclude_self_idx=None):
778        rows=[]
779        for i in range(len(df_test)):
780            input_id = df_test.loc[i, id_col]
781            pairs = [(float(dist_mat[i, j]), int(idx_mat[i, j])) for j in range(dist_mat.
               shape[1])]
782            if exclude_self_idx is not None:
783                pairs = [p for p in pairs if p[1] != exclude_self_idx(i)]
784            if len(pairs) < 2: continue
785            (d1, r1), (d2, r2) = pairs[0], pairs[1]
786            row1 = cases_df.iloc[r1]; row2 = cases_df.iloc[r2]
787            top1_id = row1.get(id_col, row1.get("Identificador", f"row_{r1}"))
788            top2_id = row2.get(id_col, row2.get("Identificador", f"row_{r2}"))
789            rows.append({"input_id": input_id,
790                        "top1_id": top1_id, "top1_dist": d1, "top1_sim%": sim_percent(d1
                        ),
791                        "top2_id": top2_id, "top2_dist": d2, "top2_sim%": sim_percent(d2
                        )})
792        out = pd.DataFrame(rows)
793        out.to_csv(out_csv_path, index=False)
794        return out
795
796    pairs_all_path = os.path.join(OUT_DIR, "knn_test_validation_pairs_all.csv")
797    pairs_all = build_pairs_csv(df_test, idx_all, dist_all, cases_all, pairs_all_path)
798    print(f"📄 {pairs_all_path}")
799
800    id_to_train_idx={}
801    if id_col in cases_train.columns:
802        for i, r in cases_train.reset_index(drop=True).iterrows():
803            id_to_train_idx[r.get(id_col, f"row_{i}")] = int(r[INDEX_KEY])
804
805    def exclude_self_idx(i):
806        _id = df_test.loc[i, id_col]
807        return id_to_train_idx.get(_id, -99999)
808
809    pairs_train_path = os.path.join(OUT_DIR, "knn_test_validation_pairs_train.csv")
810    pairs_train = build_pairs_csv(df_test, idx_train, dist_train, cases_train,
       pairs_train_path, exclude_self_idx=exclude_self_idx)
811    print(f"📄 {pairs_train_path}")
812
813    # Métricas rápidas
```

```python
814    SIM_THRESHOLD = 60.0
815    def hit_at_least_one(df_pairs, sim_thr=SIM_THRESHOLD):
816        if df_pairs.empty: return 0.0
817        return float((df_pairs[["top1_sim%","top2_sim%"]].max(axis=1) >= sim_thr).mean())*
           100.0
818
819    metrics = {
820        "built_at": datetime.now().isoformat(),
821        "p95_nn1_train": p95_nn1_train,
822        "hit@2>=60_ALL(%)": hit_at_least_one(pairs_all),
823        "hit@2>=60_TRAIN(%)": hit_at_least_one(pairs_train),
824        "space": "pre_num"
825    }
826    with open(os.path.join(OUT_DIR, "knn_validation_metrics_extended.json"), "w", encoding
       ="utf-8") as f:
827        json.dump(metrics, f, ensure_ascii=False, indent=2)
828
829    # (Opcional) por cluster si existe
830    if "Cluster_6" in df_test.columns and len(pairs_all):
831        tmp = pairs_all.copy()
832        tmp["Cluster_6"] = df_test["Cluster_6"].values[:len(tmp)]
833        agg = tmp.groupby("Cluster_6").apply(lambda g: float((g[["top1_sim%","top2_sim%"
           ]].max(axis=1) >= SIM_THRESHOLD).mean())*100.0).reset_index()
834        agg.columns = ["Cluster_6", "hit@2>=60(%)_ALL"]
835        agg.to_csv(os.path.join(OUT_DIR, "knn_test_hit3_by_cluster_cases.csv"), index=
           False)
836
837    with open(os.path.join(STEP_DIR, "README_STEP.txt"), "w", encoding="utf-8") as f:
838        f.write(f"{STEP_TAG} generado el {datetime.now().isoformat()} (BASE_DIR={BASE_DIR
           }).\n")
839
840    print("✅ Validación kNN completada (todo bajo _v2).")
841
842    """Celda — XGB N_ofertantes (desarrollo + validación + guardados producción) desde
       df_test_train_v8.xlsx"""
843
844    # ================================================================
845    # XGB N_ofertantes — desarrollo + validación + guardados producción
846    # (desde df_test_train_v8.xlsx, sin usar comparativa como fuente)
847    # ================================================================
848    import os, json, joblib
849    import numpy as np
850    import pandas as pd
851    from pathlib import Path
852    from sklearn.model_selection import train_test_split
853    from sklearn.preprocessing import LabelEncoder
854    from sklearn.metrics import accuracy_score, f1_score, recall_score, confusion_matrix
855    import xgboost as xgb
856    import matplotlib.pyplot as plt
857
858    # --- Rutas base (todo bajo _v2) ---
859    BASE_DIR  = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
860    DIR_MODELOS = BASE_DIR/"modelos"
861    DIR_VALID  = BASE_DIR/"validacion"
862    DIR_COMP   = BASE_DIR/"comparativa"   # solo para la celda de comparativa (no leemos
       aquí)
863    DIR_MODELOS.mkdir(parents=True, exist_ok=True)
864
865    # Validación: subcarpeta (mantengo el nombre que usaste en tu ejemplo)
866    vdir = DIR_VALID / "03_xgb_n_ofertantes"
867    vdir.mkdir(parents=True, exist_ok=True)
868
869    RANDOM_STATE = 42
870
871    # -------- 0) Carga df desde inputs --------
872    df_candidates = [BASE_DIR/"inputs/df_test_train_v8.xlsx", BASE_DIR/
       "df_test_train_v8.xlsx"]
873    df_path = next((p for p in df_candidates if p.exists()), None)
874    assert df_path is not None, f"No se encontró df_test_train_v8.xlsx en: {df_candidates
       }"
875    df = pd.read_excel(df_path)
876
```

```python
877      # -------- 1) Variables preferidas + alias ----------
878      XGB_VARS_PREF = [
879          "Presupuesto_licitacion_lote_c","Plazo_m","Criterio_precio_p","N_CPV",
880          "N_clasi_empresa","Mes_lici","Intervalo_lici_d_c","N_lotes",
881          "Tipo_de_Administracion_c","Tipo_de_procedimiento_c","C_juicio_valor_p_c",
882          "Provincia2"
883      ]
884      ALIASES = {
885          "Plazo_m": ["Plazo_m_c"],
886          "Criterio_precio_p": ["C_precio_p"],
887      }
888
889      def resolve_vars(preferred, df_cols, aliases):
890          used, mapping, missing = [], {}, []
891          for v in preferred:
892              cands = [v] + aliases.get(v, [])
893              found = next((c for c in cands if c in df_cols), None)
894              if found is None:
895                  missing.append(v)
896              else:
897                  used.append(found); mapping[v] = found
898          return used, mapping, missing
899
900      # Asegurar Provincia2 existe (derivamos si falta)
901      if "Provincia2" not in df.columns:
902          def _to_cp_str(x):
903              try: return str(int(x)).zfill(5)
904              except Exception:
905                  s = str(x); return s if (isinstance(s, str) and s and s.lower()!="nan")
                      else np.nan
906          if "Codigo_Postal_c" in df.columns:
907              df["Provincia2"] = df["Codigo_Postal_c"].apply(_to_cp_str).str[:2].astype(
                  "string")
908          else:
909              df["Provincia2"] = pd.Series([pd.NA]*len(df), dtype="string")
910
911      # Resolver columnas según existan en df
912      XGB_VARS_USED, VAR_MAP, MISSING_PREF = resolve_vars(XGB_VARS_PREF, set(df.columns),
         ALIASES)
913
914      # Validación de requeridos mínimos para entrenar
915      REQUIRED_MIN = {"Presupuesto_licitacion_lote_c","N_CPV","N_clasi_empresa","Mes_lici",
916                      "Intervalo_lici_d_c","N_lotes","Tipo_de_Administracion_c",
                        "Tipo_de_procedimiento_c","C_juicio_valor_p_c"}
917      missing_min = [v for v in REQUIRED_MIN if (v not in VAR_MAP and v not in XGB_VARS_USED
         )]
918      if "N_ofertantes" not in df.columns:
919          raise ValueError("✗ 'N_ofertantes' no está en df. No se puede entrenar XGB.")
920      if missing_min:
921          raise ValueError(f"✗ Faltan columnas mínimas para XGB (tras alias): {missing_min
             }")
922
923      # -------- 2) Target -> clases 0/1/2 ----------
924      def cat_n_ofertantes(n):
925          if pd.isna(n): return np.nan
926          n = int(n)
927          if n <= 4:  return 0
928          if n <= 11: return 1
929          return 2
930
931      mask = df["N_ofertantes"].notna()
932      X_raw = df.loc[mask, XGB_VARS_USED].copy()
933      y = df.loc[mask, "N_ofertantes"].apply(cat_n_ofertantes).astype(int)
934      labels_txt = ["BAJO","MEDIO","ALTO"]
935
936      # -------- 3) Split 80/20 estratificado ----------
937      X_train, X_test, y_train, y_test = train_test_split(
938          X_raw, y, test_size=0.20, random_state=RANDOM_STATE, stratify=y
939      )
940
941      # -------- 4) Encoding (LabelEncoder) en 3 categóricas ----------
942      cat_vars_pref = ["Tipo_de_Administracion_c","Tipo_de_procedimiento_c","Provincia2"]
```

```python
    # Mapeo preferida->usada para encoders
    encoders = {}
    for var_pref in cat_vars_pref:
        var_used = VAR_MAP.get(var_pref, var_pref)
        if var_used in X_train.columns:
            all_cats = pd.concat([X_train[var_used], X_test[var_used]]).astype(str).unique()
            le = LabelEncoder().fit(all_cats)
            X_train[var_used] = le.transform(X_train[var_used].astype(str))
            X_test[var_used]  = le.transform(X_test[var_used].astype(str))
            encoders[var_used] = le  # guardamos encoder con el NOMBRE REAL usado

    # Mes_lici numérica
    mes_col = VAR_MAP.get("Mes_lici", "Mes_lici")
    if mes_col in X_train.columns and not pd.api.types.is_numeric_dtype(X_train[mes_col]):
        X_train[mes_col] = pd.to_numeric(X_train[mes_col], errors="coerce").fillna(0).astype(int)
        X_test[mes_col]  = pd.to_numeric(X_test[mes_col],  errors="coerce").fillna(0).astype(int)

    # -------- 5) Modelo XGB ----------
    xgb_model = xgb.XGBClassifier(
        n_estimators=50, max_depth=4, learning_rate=0.2,
        random_state=RANDOM_STATE, eval_metric="mlogloss", verbosity=0
    )
    xgb_model.fit(X_train, y_train)

    # -------- 6) Métricas ----------
    y_pred_tr = xgb_model.predict(X_train)
    y_pred_te = xgb_model.predict(X_test)
    p_te = xgb_model.predict_proba(X_test)

    acc_tr = accuracy_score(y_train, y_pred_tr)
    acc_te = accuracy_score(y_test,  y_pred_te)
    f1_per = f1_score(y_test, y_pred_te, average=None)
    rec_per = recall_score(y_test, y_pred_te, average=None)
    f1_w   = f1_score(y_test, y_pred_te, average="weighted")

    # -------- 7) Matriz de confusión ----------
    cm = confusion_matrix(y_test, y_pred_te, labels=[0,1,2])
    plt.figure(figsize=(6,5))
    plt.imshow(cm, interpolation="nearest")
    plt.title("Matriz de Confusión — XGB N_ofertantes")
    plt.colorbar()
    ticks = np.arange(3); plt.xticks(ticks, labels_txt, rotation=45); plt.yticks(ticks, labels_txt)
    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]):
            plt.text(j, i, format(cm[i,j], "d"), ha="center", va="center")
    plt.ylabel("Reales"); plt.xlabel("Predicciones")
    plt.tight_layout(); plt.savefig(vdir/"cm_xgb.png", dpi=150); plt.close()

    # -------- 8) Tabla de métricas por clase + total ----------
    tbl = []
    for i, L in enumerate(labels_txt):
        casos_reales    = int((y_test==i).sum())
        casos_predichos = int((y_pred_te==i).sum())
        prec = (cm[i,i]/casos_predichos) if casos_predichos>0 else 0.0
        tbl.append({"Clase": L, "Casos_Reales": casos_reales, "Casos_Predichos": casos_predichos,
                    "F1_Score": float(f1_per[i]), "Recall": float(rec_per[i]), "Precision": float(prec)})
    tbl.append({"Clase": "TOTAL", "Casos_Reales": int(len(y_test)), "Casos_Predichos": int(len(y_test)),
                "F1_Score": float(f1_w), "Recall": float(acc_te), "Precision": float(acc_te)})
    pd.DataFrame(tbl).to_csv(vdir/"tabla_metricas_XGB_v8.csv", index=False)

    # -------- 9) REAL vs PREDICHO + probas ----------
    df_pred = pd.DataFrame({
        "ID_Test": np.arange(len(y_test)),
        "Real": y_test.values,
```

```python
        "Real_Label": [labels_txt[i] for i in y_test.values],
        "Predicho": y_pred_te,
        "Predicho_Label": [labels_txt[i] for i in y_pred_te],
        "Correcto": (y_test.values == y_pred_te),
        "Prob_BAJO": p_te[:,0], "Prob_MEDIO": p_te[:,1], "Prob_ALTO": p_te[:,2],
        "Prob_Max": p_te.max(axis=1)
    })
    df_pred.to_csv(vdir/"real_vs_predicho_XGB_v8.csv", index=False)

    # -------- 10) Datasets train/test codificados (traza) ----------
    df_train_final = pd.DataFrame(X_train, columns=X_train.columns)
    df_train_final["N_ofertantes_cat"]   = y_train.values
    df_train_final["N_ofertantes_label"] = ["BAJO" if i==0 else ("MEDIO" if i==1 else
    "ALTO") for i in y_train.values]
    df_test_final = pd.DataFrame(X_test, columns=X_test.columns)
    df_test_final["N_ofertantes_cat"]   = y_test.values
    df_test_final["N_ofertantes_label"] = ["BAJO" if i==0 else ("MEDIO" if i==1 else
    "ALTO") for i in y_test.values]
    df_train_final.to_csv(vdir/"df_train_XGB_v8.csv", index=False)
    df_test_final.to_csv(vdir/"df_test_XGB_v8.csv", index=False)

    # -------- 11) Importancias ----------
    imp = getattr(xgb_model, "feature_importances_", None)
    if imp is not None and len(imp) == df_train_final.shape[1]-2:
        pd.DataFrame({"Variable": df_train_final.columns[:-2], "Importancia": imp})\
            .sort_values("Importancia", ascending=False)\
            .to_csv(vdir/"importancia_variables_XGB_v8.csv", index=False)

    # -------- 12) Pack para PRODUCCIÓN ----------
    pack = {
        "modelo": xgb_model,
        "encoders": encoders,                    # dict: {col_usada: LabelEncoder}
        "variables_predictoras": XGB_VARS_USED,  # NOMBRES REALES usados (tras alias)
        "labels": ["BAJO","MEDIO","ALTO"],
        "alias_mapping": VAR_MAP,                # mapping preferida->usada (traza)
        "metricas_test": {
            "accuracy": float(acc_te),
            "f1_scores": [float(x) for x in f1_per],
            "recall_scores": [float(x) for x in rec_per],
            "f1_weighted": float(f1_w)
        },
        "version": "v8",
        "random_state": int(RANDOM_STATE)
    }
    joblib.dump(pack, DIR_MODELOS / "model_Noferta_XGB_v8.pkl")

    # -------- 13) Predict de ejemplo ----------
    ejemplo = X_test.iloc[:1].copy()
    ejemplo.to_csv(vdir/"predict_XGB_v8.csv", index=False)

    # -------- 14) Resumen + explicación ----------
    with open(vdir/"metrics.json","w",encoding="utf-8") as f:
        json.dump({"accuracy_train": float(acc_tr), "accuracy_test": float(acc_te),
                   "f1_weighted": float(f1_w),
                   "f1_per_class": dict(zip(["BAJO","MEDIO","ALTO"], [float(x) for x in
                   f1_per])),
                   "recall_per_class": dict(zip(["BAJO","MEDIO","ALTO"], [float(x) for x
                   in rec_per]))},
                  f, indent=2, ensure_ascii=False)

    with open(vdir/"explicacion_tecnica.txt","w",encoding="utf-8") as f:
        f.write(
            "XGB multiclase (BAJO/MEDIO/ALTO) para N_ofertantes con split 80/20 (rs=42). "
            "Resolución de alias (p.ej., Criterio_precio_p/C_precio_p; "
            "Plazo_m/Plazo_m_c). "
            "Encoding: LabelEncoder en Tipo_de_Administracion_c, Tipo_de_procedimiento_c,
            Provincia2 (usando nombres reales tras alias). "
            "Mes_lici numérica. Pack de producción guarda variables usadas y el mapping
            de alias."
        )

    print("✅ XGB N_ofertantes entrenado y guardado")
```

```
1072    print("→ modelos/: model_Noferta_XGB_v8.pkl")
1073    print("→ validacion/03_xgb_n_ofertantes/: metrics.json, cm_xgb.png,
        tabla_metricas_XGB_v8.csv, real_vs_predicho_XGB_v8.csv, df_train_XGB_v8.csv,
        df_test_XGB_v8.csv, predict_XGB_v8.csv, importancia_variables_XGB_v8.csv")
1074
1075    """Celda C — META del XGB (baja_comp): políticas y fallbacks"""
1076
1077    # ==========================================
1078    # META — XGB baja_comp (políticas de inputs)
1079    # ==========================================
1080    import json, joblib, numpy as np, pandas as pd
1081    from pathlib import Path
1082
1083    BASE = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
1084    DIR_MODELOS = BASE / "modelos"
1085    DIR_INPUTS  = BASE / "inputs"
1086
1087    pack_path = DIR_MODELOS / "model_Noferta_XGB_v8.pkl"
1088    assert pack_path.exists(), "Falta modelos/model_Noferta_XGB_v8.pkl"
1089    pack = joblib.load(pack_path)
1090
1091    # Cargar df de partida para medir modas (preferencia inputs/df_test_train_v8.xlsx)
1092    df_path = None
1093    for p in [DIR_INPUTS/"df_test_train_v8.xlsx", BASE/"df_test_train_v8.xlsx"]:
1094        if p.exists(): df_path = p; break
1095    assert df_path is not None, "No encuentro df_test_train_v8.xlsx para extraer modas de
        categóricas."
1096    df_all = pd.read_excel(df_path)
1097
1098    # Columnas y encoders reales del pack
1099    xgb_vars  = list(pack.get("variables_predictoras", []))
1100    xgb_encs  = pack.get("encoders", {})   # {col_usada: LabelEncoder}
1101
1102    # Mapa fallback por columna categórica = moda en df_all (si no está, usamos primer
        class_ del encoder)
1103    fallback_class_por_col = {}
1104    for col, le in xgb_encs.items():
1105        if col in df_all.columns:
1106            s = df_all[col].astype(str)
1107            if len(s.dropna()) > 0:
1108                moda = s.value_counts(dropna=True).idxmax()
1109                # Si la moda no está en el encoder (variantes), elegimos la clase más
                frecuente del encoder
1110                if moda not in set(map(str, le.classes_)):
1111                    # fallback a la clase con más ocurrencias en df_all entre las classes_
1112                    vc = {cls: int((s == str(cls)).sum()) for cls in le.classes_}
1113                    moda = max(vc, key=vc.get)
1114                fallback_class_por_col[col] = str(moda)
1115            else:
1116                fallback_class_por_col[col] = str(le.classes_[0])
1117        else:
1118            fallback_class_por_col[col] = str(le.classes_[0])
1119
1120    xgb_meta = {
1121        "artifact": "model_Noferta_XGB_v8.pkl",
1122        "variables_predictoras": xgb_vars,
1123        "encoders_columns": list(xgb_encs.keys()),
1124        "labels": list(pack.get("labels", [])),
1125        "policies": {
1126            "missing_numeric_policy": "fill_0",
1127            "cat_unknown_policy": "fallback_to_mode",  # usar moda por columna
1128            "fallback_class_por_col": fallback_class_por_col,
1129            "threshold_bajo": 0.60
1130        },
1131        "version": pack.get("version", "v8"),
1132        "random_state": int(pack.get("random_state", 42))
1133    }
1134    (DIR_MODELOS / "xgb_meta.json").write_text(json.dumps(xgb_meta, indent=2, ensure_ascii
        =False), encoding="utf-8")
1135    print("✅ Guardado modelos/xgb_meta.json")
1136
1137    """RF base (PRE mínimas) — desarrollo + validación + guardados"""
```

```python
# ============================================================
# RF base (PRE mínimas) — desarrollo + validación + guardados
# ============================================================
import os, json, joblib, numpy as np, pandas as pd, matplotlib.pyplot as plt
from pathlib import Path
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.impute import SimpleImputer
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import StratifiedKFold, RandomizedSearchCV
from sklearn.metrics import accuracy_score, f1_score, confusion_matrix

# --- Rutas base ---
BASE_DIR    = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
DIR_MODELOS = BASE_DIR / "modelos"
DIR_VALID   = BASE_DIR / "validacion"
DIR_MODELOS.mkdir(parents=True, exist_ok=True)

vdir = DIR_VALID / "04_rf_base"
vdir.mkdir(parents=True, exist_ok=True)

RANDOM_STATE = 42

# --- Cargar df y splits TRAIN/TEST ---
df_candidates = [BASE_DIR/"inputs/df_test_train_v8.xlsx", BASE_DIR/
"df_test_train_v8.xlsx"]
df_path = next((p for p in df_candidates if p.exists()), None)
assert df_path is not None, f"No se encontró df_test_train_v8.xlsx en: {df_candidates
}"
df = pd.read_excel(df_path)

assert "Cluster_6" in df.columns, "✗ Falta 'Cluster_6'."
assert "Dataset"   in df.columns, "✗ Falta 'Dataset' (TRAIN/TEST)."
df["Cluster_6"] = df["Cluster_6"].astype(str)
df["Dataset"]   = df["Dataset"].astype(str).str.upper()

df_train = df[df["Dataset"] == "TRAIN"].copy()
df_test  = df[df["Dataset"] == "TEST"].copy()

# 1) Variables PRE mínimas acordadas (sin Provincia2) + alias
MIN_PRE_PREF = [
    "Presupuesto_licitacion_lote_c","Plazo_m","Criterio_precio_p","C_precio_p",
    "N_CPV","N_clasi_empresa","Mes_lici","Intervalo_lici_d_c","N_lotes",
    "Tipo_de_Administracion_c","Tipo_de_procedimiento_c","C_juicio_valor_p_c"
]
ALIASES = {
    "Criterio_precio_p": ["C_precio_p"],  # preferimos Criterio_precio_p; si no,
    C_precio_p
    "Plazo_m": ["Plazo_m_c"],             # preferimos Plazo_m; si no, Plazo_m_c
}

def resolve_vars(preferred, df_cols, aliases):
    used = []
    mapping = {}
    # Regla especial para la variable de precio: elegir UNA (preferida si existe)
    precio_chosen = None
    for v in preferred:
        cands = [v] + aliases.get(v, [])
        found = next((c for c in cands if c in df_cols), None)
        if found is None:
            continue
        if v in ["Criterio_precio_p", "C_precio_p"]:
            if precio_chosen is not None:
                # ya elegimos una de las dos, saltamos la otra
                continue
            # si existe Criterio_precio_p, tomamos esa y marcamos; si no, tomamos
            C_precio_p
            if "Criterio_precio_p" in df_cols:
                precio_chosen = "Criterio_precio_p"
                used.append("Criterio_precio_p")
```

```python
                        mapping["Criterio_precio_p"] = "Criterio_precio_p"
                    elif "C_precio_p" in df_cols:
                        precio_chosen = "C_precio_p"
                        used.append("C_precio_p")
                        mapping["Criterio_precio_p"] = "C_precio_p"  # mapping a la usada real
                    continue
                # Plazo_m vs Plazo_m_c: elegimos una
                if v == "Plazo_m":
                    if "Plazo_m" in df_cols:
                        used.append("Plazo_m"); mapping["Plazo_m"] = "Plazo_m"
                    elif "Plazo_m_c" in df_cols:
                        used.append("Plazo_m_c"); mapping["Plazo_m"] = "Plazo_m_c"
                    continue
                used.append(found); mapping[v] = found
        # quitar duplicados conservando orden
        seen=set(); used_uniq=[]
        for c in used:
            if c not in seen:
                seen.add(c); used_uniq.append(c)
        return used_uniq, mapping

    MIN_PRE_USED, VAR_MAP = resolve_vars(MIN_PRE_PREF, set(df.columns), ALIASES)
    assert len(MIN_PRE_USED) > 0, "✗ No se resolvieron variables PRE mínimas."

    # 2) Preparar X/y (TRAIN/TEST) — asegurar Mes_lici categórica para RF
    Xtr = df_train[MIN_PRE_USED].copy()
    Xte = df_test[MIN_PRE_USED].copy()
    ytr = df_train["Cluster_6"].astype(str).copy()
    yte = df_test["Cluster_6"].astype(str).copy()

    mes_col = VAR_MAP.get("Mes_lici", "Mes_lici")
    if mes_col in Xtr.columns:
        Xtr[mes_col] = Xtr[mes_col].astype("Int64").astype("category")
        Xte[mes_col] = Xte[mes_col].astype("Int64").astype("category")

    # 3) Preprocesamiento + modelo
    try:
        oh = OneHotEncoder(handle_unknown="ignore", sparse_output=False)
    except TypeError:
        oh = OneHotEncoder(handle_unknown="ignore", sparse=False)

    num_cols = Xtr.select_dtypes(include=[np.number]).columns.tolist()
    cat_cols = [c for c in Xtr.columns if c not in num_cols]

    prep = ColumnTransformer([
        ("num", Pipeline([("imp", SimpleImputer(strategy="median")), ("sc", StandardScaler
        ())]), num_cols),
        ("cat", Pipeline([("imp", SimpleImputer(strategy="most_frequent")), ("oh", oh)]),
        cat_cols),
    ])

    rf = RandomForestClassifier(random_state=RANDOM_STATE, n_jobs=-1)
    pipe = Pipeline([("prep", prep), ("clf", rf)])

    # 4) Tuning ligero
    param_grid = {
        "clf__n_estimators": [600, 900, 1100],
        "clf__max_depth": [None, 18, 24],
        "clf__min_samples_split": [2, 5, 10],
        "clf__min_samples_leaf": [1, 2, 4],
        "clf__max_features": ["sqrt", 0.5],
        "clf__class_weight": [None, "balanced_subsample"]
    }
    cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=RANDOM_STATE)
    rs = RandomizedSearchCV(
        pipe, param_distributions=param_grid, n_iter=28,
        cv=cv, scoring="f1_macro", random_state=RANDOM_STATE, n_jobs=-1, verbose=0
    )
    rs.fit(Xtr, ytr)
    rf_base = rs.best_estimator_

    # 5) Predicciones y métricas
```

```python
1276    pred = rf_base.predict(Xte)
1277    proba = rf_base.predict_proba(Xte)
1278    acc = accuracy_score(yte, pred)
1279    f1m = f1_score(yte, pred, average="macro")
1280
1281    def per_class_acc(y_true, y_pred):
1282        out={}
1283        for cl in sorted(pd.unique(y_true), key=lambda s:int(s)):
1284            m = (y_true==cl)
1285            out[cl] = float((y_pred[m]==y_true[m]).mean()) if m.any() else np.nan
1286        return out
1287
1288    acc_c = per_class_acc(yte, pred)
1289
1290    # 6) Guardados
1291    joblib.dump(rf_base, DIR_MODELOS / "rf_base.pkl")
1292    pd.DataFrame({"y_true": yte, "y_pred": pred}).to_csv(vdir / "pred_test_rf_base.csv",
        index=False)
1293    pd.Series(acc_c).to_csv(vdir / "accuracy_por_clase.csv", header=["accuracy"])
1294    with open(vdir / "summary.json","w",encoding="utf-8") as f:
1295        json.dump({
1296            "accuracy": float(acc),
1297            "macro_f1": float(f1m),
1298            "features_used": MIN_PRE_USED,
1299            "alias_mapping": VAR_MAP,
1300            "best_params": rs.best_params_
1301        }, f, indent=2, ensure_ascii=False)
1302
1303    # Matriz de confusión
1304    labs = sorted(pd.unique(pd.concat([ytr, yte])), key=lambda s:int(s))
1305    cm = confusion_matrix(yte, pred, labels=labs)
1306    plt.figure(figsize=(6,5))
1307    plt.imshow(cm, interpolation="nearest")
1308    plt.title("Matriz de Confusión — RF base")
1309    plt.colorbar()
1310    ticks=np.arange(len(labs)); plt.xticks(ticks,labs,rotation=45, ha="right"); plt.yticks
        (ticks,labs)
1311    for i in range(cm.shape[0]):
1312        for j in range(cm.shape[1]):
1313            plt.text(j,i,format(cm[i,j],'d'),ha="center",va="center")
1314    plt.ylabel("Reales"); plt.xlabel("Predicciones")
1315    plt.tight_layout(); plt.savefig(vdir / "cm_rf_base.png", dpi=150); plt.close()
1316
1317    with open(vdir / "explicacion_tecnica.txt","w",encoding="utf-8") as f:
1318        f.write(
1319            "RF base entrenado SOLO con variables PRE mínimas (sin Provincia2). "
1320            "Pipeline: imputación (mediana/moda), OHE(handle_unknown='ignore'), escalado
                numérico (var≈1). "
1321            "Tuning ligero con RandomizedSearchCV (f1_macro, rs=42). Artefacto final en
                /modelos/rf_base.pkl."
1322        )
1323
1324    print("✅ RF base entrenado y guardado en modelos/rf_base.pkl")
1325    print(f"[TEST] Acc={acc:.4f} | Macro-F1={f1m:.4f} | Vars usadas={len(MIN_PRE_USED)}")
1326
1327    """Celda D — META del RF base (features crudas exactas + dtypes esperados)"""
1328
1329    # =========================================
1330    # META — RF base (contrato de features)
1331    # =========================================
1332    import json, joblib, sklearn
1333    from pathlib import Path
1334
1335    BASE = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
1336    DIR_MODELOS = BASE / "modelos"
1337
1338    rf_path = DIR_MODELOS / "rf_base.pkl"
1339    assert rf_path.exists(), "Falta modelos/rf_base.pkl"
1340    rf_base = joblib.load(rf_path)
1341
1342    # Extraer columnas crudas esperadas por el ColumnTransformer
1343    def _cols_from_pipe(pipe):
```

```python
        try:
            trfs = pipe.named_steps["prep"].transformers_
            cols = {}
            for name, _, c in pipe.named_steps["prep"].transformers_:
                cols[name] = list(c) if isinstance(c, list) else []
            return cols
        except Exception:
            return {"num": [], "cat": []}

    cols = _cols_from_pipe(rf_base)
    alias_mapping = globals().get("VAR_MAP", {})  # si existe en el notebook, lo
    aprovechamos

    rf_meta = {
        "artifact": "rf_base.pkl",
        "expected_raw_features": cols,
        "alias_mapping_used_in_train": alias_mapping,
        "dtypes_expect": {"Mes_lici": "category"},
        "sklearn_version": sklearn.__version__
    }
    (DIR_MODELOS / "rf_base_meta.json").write_text(json.dumps(rf_meta, indent=2,
    ensure_ascii=False), encoding="utf-8")
    print("☑ Guardado modelos/rf_base_meta.json")

    """RF_124_BASE+BC (train) + Ensemble "seguro" BASE+BC (test)"""

    # =====================================================================
    # Celda única — RF_124_BASE+BC (train) + Ensemble "seguro" BASE+BC (test)
    #    - Entrena especialista 1/2/4 con MIN_PRE + baja_comp (XGB) duplicada (k=5)
    #    - Construye ensemble seguro vs RF base y genera la gráfica final
    #    - Guarda artefactos de PRODUCCIÓN (rf_124_basebc.pkl + ensemble_meta.json)
    #    - Guarda validaciones (CSV/PNG/JSON)
    # =====================================================================
    import os, json, joblib, numpy as np, pandas as pd, matplotlib.pyplot as plt
    from pathlib import Path
    from sklearn.compose import ColumnTransformer
    from sklearn.pipeline import Pipeline
    from sklearn.preprocessing import OneHotEncoder, StandardScaler
    from sklearn.impute import SimpleImputer
    from sklearn.ensemble import RandomForestClassifier
    from sklearn.model_selection import StratifiedKFold, RandomizedSearchCV
    from sklearn.metrics import accuracy_score, f1_score, confusion_matrix

    # ---------------- Configuración -----------------
    BASE_DIR    = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
    DIR_MODELOS = BASE_DIR / "modelos"
    DIR_VALID   = BASE_DIR / "validacion"
    DIR_MODELOS.mkdir(parents=True, exist_ok=True)
    (DIR_VALID / "05_rf_124_basebc").mkdir(parents=True, exist_ok=True)
    (DIR_VALID / "06_ensemble_basebc").mkdir(parents=True, exist_ok=True)

    RANDOM_STATE   = 42
    CLASSES_124    = {"1","2","4"}
    BAJA_COMP_DUP  = 5          # ← regla acordada (k=5)
    TH_BAJO        = 0.55
    ALLOWED_DROP_1 = 0.05       # ← máx 5 pp
    REQUIRE_IMPROVE_4 = True
    REQUIRE_IMPROVE_GLOBAL = True
    MARGIN_MIN     = 0.03

    vdir_model  = DIR_VALID / "05_rf_124_basebc"
    vdir_ens    = DIR_VALID / "06_ensemble_basebc"

    # ---------------- 0) Cargar datos y prerequisitos ----------------
    df_candidates = [BASE_DIR/"inputs/df_test_train_v8.xlsx", BASE_DIR/
    "df_test_train_v8.xlsx"]
    df_path = next((p for p in df_candidates if p.exists()), None)
    assert df_path is not None, f"No se encontró df_test_train_v8.xlsx en: {df_candidates
    }"
    df = pd.read_excel(df_path)

    assert "Cluster_6" in df.columns, "✗ Falta 'Cluster_6'."
```

```
1412    assert "Dataset"   in df.columns, "✗ Falta 'Dataset' (TRAIN/TEST)."
1413    df["Cluster_6"] = df["Cluster_6"].astype(str)
1414    df["Dataset"]   = df["Dataset"].astype(str).str.upper()
1415    df_train = df[df["Dataset"] == "TRAIN"].copy()
1416    df_test  = df[df["Dataset"] == "TEST"].copy()
1417
1418    rf_base_path   = DIR_MODELOS / "rf_base.pkl"
1419    xgb_pack_path  = DIR_MODELOS / "model_Noferta_XGB_v8.pkl"
1420    assert rf_base_path.exists(), f"✗ Falta {rf_base_path}. Entrena antes RF base."
1421    assert xgb_pack_path.exists(), f"✗ Falta {xgb_pack_path}. Entrena antes XGB
        N_ofertantes."
1422    rf_base  = joblib.load(rf_base_path)
1423    xgb_pack = joblib.load(xgb_pack_path)
1424
1425    # ---------------- 1) Utilidades -----------------
1426    MIN_PRE_PREF = [
1427        "Presupuesto_licitacion_lote_c","Plazo_m","Criterio_precio_p","C_precio_p",
1428        "N_CPV","N_clasi_empresa","Mes_lici","Intervalo_lici_d_c","N_lotes",
1429        "Tipo_de_Administracion_c","Tipo_de_procedimiento_c","C_juicio_valor_p_c"
1430    ]
1431    ALIASES = {"Criterio_precio_p":["C_precio_p"], "Plazo_m":["Plazo_m_c"]}
1432
1433    def resolve_vars(preferred, df_cols, aliases):
1434        used, mapping = [], {}
1435        precio_chosen = None
1436        for v in preferred:
1437            cands = [v] + aliases.get(v, [])
1438            found = next((c for c in cands if c in df_cols), None)
1439            if found is None:
1440                continue
1441            if v in ["Criterio_precio_p", "C_precio_p"]:
1442                if precio_chosen is not None:
1443                    continue
1444                if "Criterio_precio_p" in df_cols:
1445                    precio_chosen = "Criterio_precio_p"
1446                    used.append("Criterio_precio_p"); mapping["Criterio_precio_p"] =
                        "Criterio_precio_p"
1447                elif "C_precio_p" in df_cols:
1448                    precio_chosen = "C_precio_p"
1449                    used.append("C_precio_p"); mapping["Criterio_precio_p"] = "C_precio_p"
1450                continue
1451            if v == "Plazo_m":
1452                if "Plazo_m" in df_cols:
1453                    used.append("Plazo_m"); mapping["Plazo_m"] = "Plazo_m"; continue
1454                if "Plazo_m_c" in df_cols:
1455                    used.append("Plazo_m_c"); mapping["Plazo_m"] = "Plazo_m_c"; continue
1456            used.append(found); mapping[v] = found
1457        # único conservando orden
1458        seen=set(); used_uniq=[]
1459        for c in used:
1460            if c not in seen:
1461                seen.add(c); used_uniq.append(c)
1462        return used_uniq, mapping
1463
1464    MIN_PRE_USED, VAR_MAP = resolve_vars(MIN_PRE_PREF, set(df_train.columns), ALIASES)
1465
1466    # Provincia2 para XGB si falta
1467    def _to_cp_str(x):
1468        try: return str(int(x)).zfill(5)
1469        except Exception:
1470            s = str(x); return s if (isinstance(s,str) and s and s.lower()!="nan") else np
                .nan
1471    if "Provincia2" not in df.columns:
1472        if "Codigo_Postal_c" in df.columns:
1473            df["Provincia2"] = df["Codigo_Postal_c"].apply(_to_cp_str).str[:2].astype(
                "string")
1474        else:
1475            df["Provincia2"] = pd.Series([pd.NA]*len(df), dtype="string")
1476    # (recalcular train/test por si añadimos la columna)
1477    df_train = df[df["Dataset"] == "TRAIN"].copy()
1478    df_test  = df[df["Dataset"] == "TEST"].copy()
1479
```

```python
1480    # Pack XGB → predictor de baja_comp (BAJO/NO BAJO con umbral 0.55)
1481    xgb_model  = xgb_pack["modelo"]
1482    xgb_vars   = xgb_pack["variables_predictoras"]   # nombres REALES usados en el pack
1483    xgb_encs   = xgb_pack["encoders"]                # {col_usada: LabelEncoder}
1484    xgb_labels = list(xgb_pack["labels"])
1485    idx_bajo   = xgb_labels.index("BAJO") if "BAJO" in xgb_labels else 0
1486
1487    def comp_predict_bajo(df_in: pd.DataFrame) -> pd.Series:
1488        Xp = pd.DataFrame(index=df_in.index)
1489        for col in xgb_vars:
1490            Xp[col] = df_in[col] if col in df_in.columns else 0
1491        for col, le in xgb_encs.items():
1492            if col in Xp.columns:
1493                Xp[col] = le.transform(Xp[col].astype(str))
1494            else:
1495                Xp[col] = 0
1496        if "Mes_lici" in Xp.columns and not pd.api.types.is_numeric_dtype(Xp["Mes_lici"]):
1497            Xp["Mes_lici"] = pd.to_numeric(Xp["Mes_lici"], errors="coerce").fillna(0).
                astype(int)
1498        p = xgb_model.predict_proba(Xp[xgb_vars])[:, idx_bajo]
1499        return pd.Series((p >= TH_BAJO).astype(np.int8), index=df_in.index, name=
            "baja_comp")
1500
1501    # Dataset BASE+BC (MIN_PRE + baja_comp categórica duplicada k veces)
1502    def build_basebc(df_in: pd.DataFrame) -> pd.DataFrame:
1503        X = df_in[MIN_PRE_USED].copy()
1504        mes_col = VAR_MAP.get("Mes_lici","Mes_lici")
1505        if mes_col in X.columns:
1506            X[mes_col] = X[mes_col].astype("Int64").astype("category")
1507        bc = comp_predict_bajo(df_in).astype("category")
1508        X["baja_comp"] = bc
1509        for k in range(1, BAJA_COMP_DUP+1):
1510            X[f"baja_comp_dup{k}"] = X["baja_comp"]
1511        return X
1512
1513    # ---------------- 2) Preparar datos especialista 1/2/4 ----------------
1514    Xtr_basebc = build_basebc(df_train)
1515    Xte_basebc = build_basebc(df_test)
1516    pd.concat([df_train[["Cluster_6"]], Xtr_basebc], axis=1).to_csv(vdir_model/
        "train_BASEBC_124.csv", index=False)
1517    pd.concat([df_test[["Cluster_6"]],  Xte_basebc],  axis=1).to_csv(vdir_model/
        "test_BASEBC_124.csv",  index=False)
1518
1519    mask_tr = df_train["Cluster_6"].isin(CLASSES_124)
1520    X124_tr = Xtr_basebc.loc[mask_tr].copy()
1521    y124_tr = df_train.loc[mask_tr, "Cluster_6"].astype(str).copy()
1522
1523    # ---------------- 3) Entrenar RF_124_BASE+BC ----------------
1524    num_cols = X124_tr.select_dtypes(include=[np.number]).columns.tolist()
1525    cat_cols = [c for c in X124_tr.columns if c not in num_cols]
1526    try:
1527        oh = OneHotEncoder(handle_unknown="ignore", sparse_output=False)
1528    except TypeError:
1529        oh = OneHotEncoder(handle_unknown="ignore", sparse=False)
1530
1531    prep = ColumnTransformer([
1532        ("num", Pipeline([("imp", SimpleImputer(strategy="median")), ("sc", StandardScaler
            ())]), num_cols),
1533        ("cat", Pipeline([("imp", SimpleImputer(strategy="most_frequent")), ("oh", oh)]),
            cat_cols),
1534    ])
1535    rf = RandomForestClassifier(random_state=RANDOM_STATE, n_jobs=-1)
1536    pipe = Pipeline([("prep", prep), ("clf", rf)])
1537
1538    param_grid = {
1539        "clf__n_estimators":[600,900,1100],
1540        "clf__max_depth":[None,18,24],
1541        "clf__min_samples_split":[2,5,10],
1542        "clf__min_samples_leaf":[1,2,4],
1543        "clf__max_features":["sqrt",0.5],
1544        "clf__class_weight":[None,"balanced_subsample"]
1545    }
```

```python
1546    cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=RANDOM_STATE)
1547    rs = RandomizedSearchCV(pipe, param_distributions=param_grid, n_iter=28, cv=cv,
1548                            scoring="f1_macro", random_state=RANDOM_STATE, n_jobs=-1,
                                verbose=0)
1549    rs.fit(X124_tr, y124_tr)
1550    rf_124_basebc = rs.best_estimator_
1551    joblib.dump(rf_124_basebc, DIR_MODELOS / "rf_124_basebc.pkl")
1552    print("📄 Guardado: modelos/rf_124_basebc.pkl")
1553
1554    # ---------------- 4) Predicciones base y especialista en TEST ----------------
1555    # A) RF base
1556    def required_cols_from_pipe(pipe):
1557        try:
1558            trfs = pipe.named_steps["prep"].transformers_
1559            cols = []
1560            for _,_, cols_i in trfs:
1561                if isinstance(cols_i, list): cols.extend(cols_i)
1562            # único preservando orden
1563            out = []
1564            for c in cols:
1565                if c not in out: out.append(c)
1566            return out
1567        except Exception:
1568            return None
1569
1570    req_base_cols = required_cols_from_pipe(rf_base) or [c for c in df_test.columns if c
        != "Cluster_6"]
1571    Xte_base = df_test[[c for c in req_base_cols if c in df_test.columns]].copy()
1572    if "Mes_lici" in Xte_base.columns:
1573        Xte_base["Mes_lici"] = Xte_base["Mes_lici"].astype("Int64").astype("category")
1574
1575    proba_base = rf_base.predict_proba(Xte_base)
1576    labels_base = [str(c) for c in rf_base.named_steps["clf"].classes_]
1577    pred_base = np.array([labels_base[i] for i in np.argmax(proba_base, axis=1)])
1578
1579    # B) Especialista BASE+BC
1580    proba_bb = rf_124_basebc.predict_proba(Xte_basebc)
1581    labels_bb = [str(c) for c in rf_124_basebc.named_steps["clf"].classes_]
1582    pred_bb = np.array([labels_bb[i] for i in np.argmax(proba_bb, axis=1)])
1583
1584    y_true = df_test["Cluster_6"].astype(str).to_numpy()
1585
1586    # ---------------- 5) Ensemble "seguro" (flips iterativos con margen)
        ----------------
1587    def safe_ensemble(pred_b, proba_b, labels_b, pred_124, proba_124, labels_124):
1588        pred = pred_b.copy()
1589        scope = np.isin(pred_b, list(CLASSES_124))
1590        idx_b = {lb:i for i,lb in enumerate(labels_b)}
1591        idx_s = {lb:i for i,lb in enumerate(labels_124)}
1592        cands = []
1593        for i in np.where(scope)[0]:
1594            base_cls = pred_b[i]
1595            alt_cls  = pred_124[i]
1596            if alt_cls != base_cls and alt_cls in CLASSES_124 and base_cls in idx_b and
            alt_cls in idx_s:
1597                mb = proba_b[i][idx_b[base_cls]]
1598                ma = proba_124[i][idx_s[alt_cls]]
1599                if (ma - mb) >= MARGIN_MIN:
1600                    cands.append((ma-mb, i, alt_cls))
1601        cands.sort(reverse=True)
1602
1603        def acc_cls(y, yp, cls):
1604            m=(y==cls);
1605            return (yp[m]==y[m]).mean() if m.any() else np.nan
1606
1607        acc0   = accuracy_score(y_true, pred_b)
1608        acc1_0 = acc_cls(y_true, pred_b, "1")
1609        acc4_0 = acc_cls(y_true, pred_b, "4")
1610
1611        flipped = []
1612        cur = pred.copy()
1613        for _, i, alt in cands:
```

```python
                tmp = cur.copy()
                tmp[i] = alt
                acc   = accuracy_score(y_true, tmp)
                acc1  = acc_cls(y_true, tmp, "1")
                acc4  = acc_cls(y_true, tmp, "4")
                condG = (not REQUIRE_IMPROVE_GLOBAL) or (acc >= acc0 - 1e-12)
                cond4 = (not REQUIRE_IMPROVE_4)      or (acc4 >= acc4_0 - 1e-12)
                cond1 = (np.isnan(acc1_0) or np.isnan(acc1)) or (acc1 >= acc1_0 -
                    ALLOWED_DROP_1 - 1e-12)
                if condG and cond4 and cond1:
                    cur = tmp
                    acc0, acc1_0, acc4_0 = acc, acc1, acc4
                    flipped.append(i)
        return cur, flipped

    ens_bb, flips_bb = safe_ensemble(pred_base, proba_base, labels_base, pred_bb, proba_bb
    , labels_bb)

    # ---------------- 6) Métricas + gráfica final ----------------
    def per_class_acc(y, yp):
        d={}
        for lb in sorted(pd.unique(y), key=lambda s:int(s)):
            m=(y==lb); d[str(lb)] = (yp[m]==y[m]).mean() if m.any() else np.nan
        return pd.Series(d)

    acc_b  = accuracy_score(y_true, pred_base);  f1_b  = f1_score(y_true, pred_base,
    average="macro")
    acc_bb = accuracy_score(y_true, ens_bb);     f1_bb = f1_score(y_true, ens_bb,
    average="macro")

    pa_b  = per_class_acc(y_true, pred_base)
    pa_bb = per_class_acc(y_true, ens_bb)

    print(f"[TEST] RF base          -> Acc={acc_b:.4f} | Macro-F1={f1_b:.4f}")
    print(f"[TEST] Ensemble BASE+BC  -> Acc={acc_bb:.4f} | Macro-F1={f1_bb:.4f} | flips={
    len(flips_bb)}")

    # Gráfica final (barras por clase: base vs ensemble BASE+BC)
    labs = sorted(pa_b.index, key=lambda s:int(s))
    x = np.arange(len(labs)); w = 0.35
    plt.figure(figsize=(11,4))
    plt.bar(x - w/2, [pa_b[l]  for l in labs], width=w, label="RF base")
    plt.bar(x + w/2, [pa_bb[l] for l in labs], width=w, label="Ensemble BASE+BC")
    for i,l in enumerate(labs):
        for off,val in [(-w/2,pa_b[l]),(w/2,pa_bb[l])]:
            plt.text(x[i]+off, val+0.01, f"{val:.2f}", ha="center", va="bottom", fontsize=
                9)
    plt.xticks(x, labs); plt.ylabel("Accuracy"); plt.title("Accuracy por clase — RF base
    vs Ensemble BASE+BC")
    plt.legend(); plt.tight_layout(); plt.savefig(vdir_ens /
    "acc_por_clase_base_vs_ensemble.png", dpi=160); plt.close()

    # Matriz de confusión del ensemble
    cm = confusion_matrix(y_true, ens_bb, labels=labs)
    plt.figure(figsize=(7,6)); plt.imshow(cm, interpolation='nearest'); plt.title("CM —
    Ensemble BASE+BC (1/2/4)"); plt.colorbar()
    ticks=np.arange(len(labs)); plt.xticks(ticks, labs, rotation=45, ha='right'); plt.
    yticks(ticks, labs)
    for i in range(cm.shape[0]):
        for j in range(cm.shape[1]): plt.text(j,i,format(cm[i,j],'d'),ha="center",va=
            "center")
    plt.ylabel("Reales"); plt.xlabel("Predicciones")
    plt.tight_layout(); plt.savefig(vdir_ens / "cm_ensemble_BASEBC.png", dpi=160); plt.
    close()

    # ---------------- 7) Guardados (CSV/JSON) ----------------
    pd.DataFrame({
        "y_true": y_true,
        "pred_rf_base": pred_base,
        "pred_ens_bb":  ens_bb,
        "changed_bb":   (pred_base!=ens_bb).astype(int)
    }).to_csv(vdir_ens / "predicciones_comparadas_test.csv", index=False)
```

```python
1674
1675    pa_b.to_csv(vdir_ens / "per_class_rf_base.csv")
1676    pa_bb.to_csv(vdir_ens / "per_class_ensemble_BASEBC.csv")
1677
1678    with open(DIR_MODELOS / "ensemble_meta.json","w",encoding="utf-8") as f:
1679        json.dump({
1680            "strategy": "rf_base_vs_rf124_basebc_flips_iterativos",
1681            "security_rules": {
1682                "global_not_worse": True,
1683                "class4_not_worse": True,
1684                "class1_drop_pp_max": 5
1685            },
1686            "margin_min": MARGIN_MIN,
1687            "uses": {
1688                "rf_base": "rf_base.pkl",
1689                "rf_124_basebc": "rf_124_basebc.pkl",
1690                "xgb_pack": "model_Noferta_XGB_v8.pkl"
1691            },
1692            "baja_comp": {"threshold": TH_BAJO, "duplicated_times": BAJA_COMP_DUP}
1693        }, f, indent=2, ensure_ascii=False)
1694
1695    with open(vdir_ens / "summary_ensemble.json","w",encoding="utf-8") as f:
1696        json.dump({
1697            "acc_rf_base": float(acc_b), "macro_f1_rf_base": float(f1_b),
1698            "acc_ensemble_BASEBC": float(acc_bb), "macro_f1_ensemble_BASEBC": float(f1_bb
1699            ),
1699            "flips_BASEBC": int(len(flips_bb)),
1700            "ALLOWED_DROP_1": ALLOWED_DROP_1,
1701            "REQUIRE_IMPROVE_4": REQUIRE_IMPROVE_4,
1702            "REQUIRE_IMPROVE_GLOBAL": REQUIRE_IMPROVE_GLOBAL,
1703            "MARGIN_MIN": MARGIN_MIN,
1704            "baja_comp_dup": BAJA_COMP_DUP,
1705            "baja_comp_threshold": TH_BAJO
1706        }, f, indent=2, ensure_ascii=False)
1707
1708    with open(vdir_ens / "explicacion_tecnica.txt","w",encoding="utf-8") as f:
1709        f.write(
1710            "RF_124_BASE+BC entrenado con MIN_PRE + 'baja_comp' (XGB) duplicada k=5. "
1711            "Ensemble 'seguro' vs RF base con flips iterativos por margen (≥0.03) y
1711            reglas de back-off: "
1712            "no empeorar accuracy global ni clase 4; clase 1 no cae >5 pp. "
1713            "Artefactos de producción en /modelos y validaciones en
1713            /validacion/06_ensemble_basebc."
1714        )
1715
1716    print("✅ Especialista y Ensemble BASE+BC listos. Artefactos en modelos/ y
1716    validacion/06_ensemble_basebc/")
1717
1718    """Celda E — META del RF_124_BASE+BC (MIN_PRE_USED + duplicadas BC exactas)"""
1719
1720    # ==========================================
1721    # META — RF_124_BASE+BC (contrato de features)
1722    # ==========================================
1723    import json, joblib, sklearn
1724    from pathlib import Path
1725
1726    BASE = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
1727    DIR_MODELOS = BASE / "modelos"
1728
1729    rf124_path = DIR_MODELOS / "rf_124_basebc.pkl"
1730    assert rf124_path.exists(), "Falta modelos/rf_124_basebc.pkl"
1731    rf_124 = joblib.load(rf124_path)
1732
1733    def _cols_from_pipe(pipe):
1734        try:
1735            trfs = pipe.named_steps["prep"].transformers_
1736            cols = {}
1737            for name, _, c in pipe.named_steps["prep"].transformers_:
1738                cols[name] = list(c) if isinstance(c, list) else []
1739            return cols
1740        except Exception:
1741            return {"num": [], "cat": []}
```

```python
cols = _cols_from_pipe(rf_124)

# Detectar duplicadas de baja_comp exactamente como las vio el modelo
dup_names = [c for c in (cols.get("cat", []) + cols.get("num", [])) if c.startswith(
    "baja_comp")]
dup_k = max([int(c.replace("baja_comp_dup","")) for c in dup_names if c != "baja_comp"
    ] or [0])

alias_mapping = globals().get("VAR_MAP", {})  # si existe
rf124_meta = {
    "artifact": "rf_124_basebc.pkl",
    "expected_raw_features": cols,
    "alias_mapping_used_in_train": alias_mapping,
    "baja_comp_dup_k": dup_k,
    "baja_comp_dup_names": [c for c in dup_names if c != "baja_comp"],
    "requires_xgb_pack": "model_Noferta_XGB_v8.pkl",
    "sklearn_version": sklearn.__version__
}
(DIR_MODELOS / "rf_124_basebc_meta.json").write_text(json.dumps(rf124_meta, indent=2,
    ensure_ascii=False), encoding="utf-8")
print("✅ Guardado modelos/rf_124_basebc_meta.json")

"""Celda F — META del Ensemble (verificación)"""

# ==========================================
# META — Ensemble (verificación y dependencias)
# ==========================================
import json
from pathlib import Path

BASE = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
DIR_MODELOS = BASE / "modelos"

ens_path = DIR_MODELOS / "ensemble_meta.json"
assert ens_path.exists(), "Falta modelos/ensemble_meta.json"
ens = json.loads(ens_path.read_text(encoding="utf-8"))

ens.setdefault("depends_on", {})
ens["depends_on"].update({
    "rf_base": "rf_base.pkl",
    "rf_124_basebc": "rf_124_basebc.pkl",
    "xgb_pack": "model_Noferta_XGB_v8.pkl",
    "preprocess": "preprocess_pipeline.pkl",
    "centroids": "centroides_pre.json"
})

ens_path.write_text(json.dumps(ens, indent=2, ensure_ascii=False), encoding="utf-8")
print("✅ Ensemble meta verificado/actualizado (modelos/ensemble_meta.json)")

"""Celda G — Resumen rápido de contratos"""

# ==========================================
# META — Resumen de contratos guardados
# ==========================================
from pathlib import Path

BASE = Path("/content/drive/MyDrive/_Pipeline_desarrollo_modelos_v2")
DIR_MODELOS = BASE / "modelos"
files = [
    "preprocess_meta.json",
    "knn_params.json",
    "knn_num_pipeline.joblib",
    "xgb_meta.json",
    "rf_base_meta.json",
    "rf_124_basebc_meta.json",
    "ensemble_meta.json"
]
print("Contratos/artefactos meta en modelos/:")
for f in files:
    p = DIR_MODELOS / f
    print(" -", f, "OK" if p.exists() else "✗")
```