

```

1  ##### SCRIPT 17 - PYTHON
2  # PIPELINE EJECUCION DE MODELOS EN PYTHON
3  # PIPELINE DE PREDICCIÓN CON ARCHIVOS PKL
4  #
5  # =====
6  # -*- coding: utf-8 -*-
7
8  # =====
9  # PRODUCCIÓN (limpio) - Setup + carga de artefactos y utils
10 # =====
11 import os, json, re, unicodedata
12 import numpy as np
13 import pandas as pd
14 from pathlib import Path
15 import joblib
16
17 # Rutas
18 BASE_DIR      = Path("/content/drive/MyDrive/_Pipeline_produccion_prediccion")
19 DIR_MODELOS   = BASE_DIR / "modelos"
20 DIR_INPUTS    = BASE_DIR / "inputs"
21 DIR_OUTPUTS   = BASE_DIR / "outputs"
22 for d in [DIR_MODELOS, DIR_INPUTS, DIR_OUTPUTS]:
23     d.mkdir(parents=True, exist_ok=True)
24
25 # Artefactos requeridos
26 paths = {
27     "preprocess_meta":          DIR_MODELOS/"preprocess_meta.json",
28     "preprocess":                DIR_MODELOS/"preprocess_pipeline.pkl",
29     "centroides":               DIR_MODELOS/"centroides_pre.json",
30     "rf_base":                  DIR_MODELOS/"rf_base.pkl",
31     "rf_base_meta":             DIR_MODELOS/"rf_base_meta.json",
32     "rf_124":                   DIR_MODELOS/"rf_124_basebc.pkl",
33     "rf_124_meta":              DIR_MODELOS/"rf_124_basebc_meta.json",
34     "xgb_pack":                 DIR_MODELOS/"model_Noferta_XGB_v8.pkl",
35     "xgb_meta":                 DIR_MODELOS/"xgb_meta.json",
36     "knn_train":                DIR_MODELOS/"knn_pre_train.joblib",
37     "knn_all":                  DIR_MODELOS/"knn_pre_all.joblib",
38     "knn_params":               DIR_MODELOS/"knn_params.json",
39     "knn_num_pipe":             DIR_MODELOS/"knn_num_pipeline.joblib",
40     "knn_catalog":              DIR_MODELOS/"knn_catalog_all.parquet",
41     "knn_catalog_meta":         DIR_MODELOS/"knn_catalog_meta.json",
42     "ensemble_meta":            DIR_MODELOS/"ensemble_meta.json",
43 }
44
45 # Checks mínimos (dejamos catálogo como opcional)
46 for k, p in paths.items():
47     if k in ["knn_catalog", "knn_catalog_meta"]:
48         continue
49     assert p.exists(), f"✗ Falta {p}"
50
51 # Cargas
52 preprocess_meta = json.loads(paths["preprocess_meta"].read_text(encoding="utf-8"))
53 preprocess       = joblib.load(paths["preprocess"])
54 centroides      = json.loads(paths["centroides"].read_text(encoding="utf-8"))
55 rf_base         = joblib.load(paths["rf_base"])
56 rf_base_meta    = json.loads(paths["rf_base_meta"].read_text(encoding="utf-8"))
57 rf_124          = joblib.load(paths["rf_124"])
58 rf_124_meta     = json.loads(paths["rf_124_meta"].read_text(encoding="utf-8"))
59 xgb_pack        = joblib.load(paths["xgb_pack"])
60 xgb_meta        = json.loads(paths["xgb_meta"].read_text(encoding="utf-8"))
61 nn_train        = joblib.load(paths["knn_train"])
62 nn_all          = joblib.load(paths["knn_all"])
63 knn_params      = json.loads(paths["knn_params"].read_text(encoding="utf-8"))
64 knn_num_pipe    = joblib.load(paths["knn_num_pipe"])
65 ensemble_meta   = json.loads(paths["ensemble_meta"].read_text(encoding="utf-8"))
66
67 # Catálogo (opcional pero recomendado)
68 catalog, id_col = None, None
69 if paths["knn_catalog"].exists():
70     catalog = pd.read_parquet(paths["knn_catalog"])
71     for cand in ["Identificador", "identificador", "ID", "Id", "id_lote", "id_contrato"]:
72         if cand in catalog.columns:

```

```

73         id_col = cand
74         break
75
76 # Helpers de normalización de texto
77 def _strip_accents(s: str) -> str:
78     s = unicodedata.normalize("NFKD", s)
79     return "".join(ch for ch in s if not unicodedata.combining(ch))
80
81 def _norm_text(x):
82     if pd.isna(x): return x
83     s = str(x).strip().lower()
84     s = _strip_accents(s)
85     s = re.sub(r"\s+", " ", s)
86     return s
87
88 # Mapa meses ES/EN -> 1..12
89 MES_MAP = {
90     "1": "1", "enero": "1", "ene": "1", "january": "1", "jan": "1",
91     "2": "2", "febrero": "2", "feb": "2", "february": "2",
92     "3": "3", "marzo": "3", "mar": "3", "march": "3",
93     "4": "4", "abril": "4", "abr": "4", "april": "4", "apr": "4",
94     "5": "5", "mayo": "5", "may": "5",
95     "6": "6", "junio": "6", "jun": "6", "june": "6",
96     "7": "7", "julio": "7", "jul": "7", "july": "7",
97     "8": "8", "agosto": "8", "ago": "8", "august": "8", "aug": "8",
98     "9": "9", "septiembre": "9", "setiembre": "9", "sep": "9", "september": "9",
99     "10": "10", "octubre": "10", "oct": "10", "october": "10",
100    "11": "11", "noviembre": "11", "nov": "11", "november": "11",
101    "12": "12", "diciembre": "12", "dic": "12", "december": "12", "dec": "12"
102 }
103
104 # POST variables que se reportarán
105 POST_PREFS = {
106     "N_ofertantes": ["N_ofertantes"],
107     "Precio_adjudicacion": ["Precio_adjudicacion", "Importe_adjudicacion_lote",
108                             "Importe_adjudicacion"],
109     "Baja_p": ["Baja_p", "Baja_percent", "Baja_porcentaje", "Baja_p_c"],
110     "Modi_eco_p": ["Modi_eco_p", "Modificacion_economica_p", "Modif_eco_p"],
111     "Precio_final_e_c": ["Precio_final_e_c", "Importe_final_e_c", "Importe_final_c"],
112     "Ahorro_final": ["Ahorro_final", "Ahorro_total", "Ahorro_final_c"]
113 }
114 def resolve_post_columns(df):
115     mapping = {}
116     for pref, cands in POST_PREFS.items():
117         found = next((c for c in cands if c in df.columns), None)
118         mapping[pref] = found
119     return mapping
120 POST_MAP = resolve_post_columns(catalog) if catalog is not None else {k:None for k in POST_PREFS}
121
122 """Celda 1 - Validación/normalización común de inputs (solo PRE)"""
123
124 # =====#
125 # VALIDACIÓN (común) - aliases, Mes_lici, Provincia2, tipos
126 # =====#
127
128 def _derive_provincia2(cp):
129     if pd.isna(cp): return np.nan
130     try:
131         s = str(int(cp)).zfill(5)
132     except Exception:
133         s = re.sub(r"\D", "", str(cp)).zfill(5)
134     return s[:2]
135
136 def validate_normalize_df(df_in: pd.DataFrame, strict: bool=True):
137     df = df_in.copy()
138     msgs = []
139
140     # Aliases numéricos (preferidos)
141     if "Criterio_precio_p" not in df.columns and "C_precio_p" in df.columns:
142         df["Criterio_precio_p"] = df["C_precio_p"]
143         msgs.append({"row": "*", "field": "Criterio_precio_p", "level": "info", "code": "alias", "hint": "C_precio_p→Criterio_precio_p"})
144
145     if "Plazo_m" not in df.columns and "Plazo_m_c" in df.columns:

```

```

142 df["Plazo_m"] = df["Plazo_m_c"]
143 msgs.append({"row":"*", "field":"Plazo_m", "level":"info", "code":"alias", "hint":
144   :"Plazo_m_c→Plazo_m"})
145
146 # Provincia2 desde CP si falta
147 if "Provincia2" not in df.columns and "Codigo_Postal_c" in df.columns:
148     df["Provincia2"] = df["Codigo_Postal_c"].apply(_derive_provincia2)
149     msgs.append({"row":"*", "field": "Provincia2", "level": "info", "code": "derive",
150       "hint": "Derivado de Codigo_Postal_c"})
151
152 # Mes_lici (1..12)
153 if "Mes_lici" in df.columns:
154     old = df["Mes_lici"].copy()
155     def _map_mes(x):
156         if pd.isna(x): return np.nan
157         s = _norm_text(x)
158         return int(MES_MAP[s]) if s in MES_MAP else np.nan
159     df["Mes_lici"] = df["Mes_lici"].apply(_map_mes)
160     bad = df["Mes_lici"].isna() & ~old.isna()
161     for i in np.where(bad)[0]:
162         msgs.append({"row": int(i), "field": "Mes_lici", "level": "error" if strict
163           else "warning",
164             "code": "month_unmapped", "original": old.iloc[i], "hint":
165               "Usa 1-12 o ene/enero/..."})
166
167 else:
168     df["Mes_lici"] = np.nan
169     msgs.append({"row": "*", "field": "Mes_lici", "level": "warning", "code": "missing",
170       "hint": "Se imputará"})
171
172 # Enforce tipos numéricos básicos
173 NUM_USED = ["N_lotes", "Criterio_precio_p", "N_CPV", "N_clasi_empresa",
174   "Presupuesto_licitacion_lote_c", "Plazo_m", "C_juicio_valor_p_c",
175   "Intervalo_lici_d_c"]
176 for col in NUM_USED:
177     if col not in df.columns: df[col] = np.nan
178     df[col] = pd.to_numeric(df[col], errors="coerce")
179
180 # Categóricas a string (se codificarán en adaptadores)
181 CAT_USED = ["Tipo_de_contrato_c", "Tipo_de_Administracion_c",
182   "Tipo_de_procedimiento_c",
183   "Tramitacion_c", "Tipo_ganador_lote_c", "Provincia2"]
184 for col in CAT_USED:
185     if col not in df.columns: df[col] = pd.NA
186     df[col] = df[col].astype("string")
187
188 # Dataframe normalizado de salida
189 KEEP_COLS = sorted(set(NUM_USED + CAT_USED + ["Mes_lici", "Provincia2",
190   "Codigo_Postal_c", "C_precio_p", "Plazo_m_c"]))
191 df_norm = df[[c for c in KEEP_COLS if c in df.columns]].copy()
192
193 # Persistencia del reporte (opcional)
194 (DIR_OUTPUTS / "report_validacion.json").write_text(json.dumps(msgs, indent=2,
195   ensure_ascii=False), encoding="utf-8")
196 return df_norm, msgs
197
198 """Celda 2 – Adaptador y predicción de XGB baja_comp (con políticas de
199 xgb_meta.json)"""
200
201 # =====
202 # XGB – baja_comp (p_bajo y flag) siguiendo xgb_meta.json
203 # =====
204
205 from sklearn.preprocessing import LabelEncoder
206
207 def predict_baja_comp_xgb(df_norm: pd.DataFrame):
208     meta = xgb_meta["policies"]
209     xgb_model = xgb_pack["modelo"]
210     xgb_vars = xgb_meta["variables_predictoras"]
211     encoders = xgb_pack.get("encoders", {})
212     labels = xgb_meta.get("labels", ["BAJO", "MEDIO", "ALTO"])
213     idx_bajo = labels.index("BAJO") if "BAJO" in labels else 0
214     th_bajo = float(meta.get("threshold_bajo", 0.38))

```

```

204 # Construcción exacta de X (en orden)
205 X = pd.DataFrame(index=df_norm.index)
206 for col in xgb_vars:
207     if col in df_norm.columns:
208         X[col] = df_norm[col]
209     else:
210         # Política de faltantes numéricos: fill_0
211         X[col] = 0
212
213 # Encoders de categóricas: fallback_to_mode por columna
214 fallback_map = meta.get("fallback_class_por_col", {})
215 for col, le in encoders.items():
216     if col in X.columns:
217         s = X[col].astype(str)
218         # Si valor no en repertorio, sustituir por fallback de la columna
219         classes = set(map(str, le.classes_))
220         fb = str(fallback_map.get(col, next(iter(classes))))
221         s = s.where(s.isin(classes), other=fb)
222         X[col] = le.transform(s)
223     else:
224         # columna prevista pero ausente -> 0
225         X[col] = 0
226
227 # Mes_lici a int si está
228 if "Mes_lici" in X.columns and not pd.api.types.is_integer_dtype(X["Mes_lici"]):
229     X["Mes_lici"] = pd.to_numeric(X["Mes_lici"], errors="coerce").fillna(0).astype(int)
230
231 proba = xgb_model.predict_proba(X[xgb_vars])
232 p_bajo = proba[:, idx_bajo]
233 baja_comp = (p_bajo >= th_bajo).astype(int)
234 return pd.Series(baja_comp, index=X.index, name="baja_comp"), pd.Series(p_bajo,
235 index=X.index, name="p_bajo")
236 """
237 Celda 3 – Adaptador y predicción RF base (usando rf_base_meta.json)
238 #
239 =====
240 # RF base – Adaptador ESTRICTO (alias de train + orden + dtypes) + diagnóstico OHE
241 # Sustituye o coloca esta celda tras la actual "Celda 3 – RF base" (redefine
242 predict_rf_base)
243 # Requiere: rf_base, rf_base_meta, DIR_OUTPUTS ya cargados (Celda 0)
244 #
245 =====
246 from typing import Tuple, List, Dict, Optional
247 import numpy as np
248 import pandas as pd
249 import json
250 from datetime import datetime
251
252 def _ensure_train_used_aliases(df_norm: pd.DataFrame, rf_meta: dict) -> Tuple[pd.
253 DataFrame, List[dict]]:
254     """
255     Garantiza que existan en df_norm las columnas EXACTAS con las que se entrenó el
256     RF base
257     (según rf_base_meta['alias_mapping_used_in_train']). Si faltan, intenta crearlas
258     desde
259     su 'preferred' o desde alias comunes (C_precio_p <-> Criterio_precio_p, Plazo_m
260     <-> Plazo_m_c).
261     """
262     msgs: List[dict] = []
263     out = df_norm.copy()
264
265     alias_train: Dict[str, str] = rf_meta.get("alias_mapping_used_in_train", {}) or {}
266     # Fallbacks bidireccionales típicos
267     fallback_alias = {
268         "C_precio_p": ["Criterio_precio_p"],
269         "Criterio_precio_p": ["C_precio_p"],
270         "Plazo_m_c": ["Plazo_m"],
271         "Plazo_m": ["Plazo_m_c"],
272     }

```

```

265
266
267     for preferred, used in alias_train.items():
268         # Si la columna "used" de train ya existe, nada que hacer
269         if used in out.columns:
270             continue
271
272         # 1) si tengo la preferida, la copio a la usada por train
273         if preferred in out.columns:
274             out[used] = out[preferred]
275             msgs.append({"level": "info", "code": "alias_materialized", "used": used, "from": "preferred"})
276             continue
277
278         # 2) si no, pruebo con fallbacks simétricos típicos
279         fb_cands = fallback_alias.get(used, []) + fallback_alias.get(preferred, [])
280         created = False
281         for alt in fb_cands:
282             if alt in out.columns:
283                 out[used] = out[alt]
284                 msgs.append({"level": "info", "code": "alias_fallback_materialized", "used": used, "from": alt})
285                 created = True
286                 break
287
288         # 3) si sigue sin estar, lo creo vacío (NaN) y lo reporto
289         if not created:
290             out[used] = np.nan
291             msgs.append({
292                 "level": "warning", "code": "alias_missing_filled_nan",
293                 "used": used, "hint": "Columna de entrenamiento no encontrada; creada como NaN"
294             })
295
296     return out, msgs
297
298 def _subset_exact_order_for_pipe(df_src: pd.DataFrame, rf_meta: dict) -> Tuple[pd.
299 DataFrame, List[dict]]:
300     """
301     Subconjunta en el ORDEN EXACTO las columnas crudas que el ColumnTransformer
302     espera.
303     Si falta alguna, la crea como NaN (y lo reporta). 'Mes_lici' se castea a category
304     si aplica.
305     Convierte numéricas a tipo numérico (el imputer de mediana está dentro del pipe).
306     """
307     msgs: List[dict] = []
308
309     exp: Dict[str, List[str]] = rf_meta.get("expected_raw_features", {}) or {}
310     cols_num = list(exp.get("num", []))
311     cols_cat = list(exp.get("cat", []))
312     ordered = cols_num + cols_cat
313
314     # Construir frame con orden exacto
315     X = pd.DataFrame(index=df_src.index, columns=ordered, dtype="object")
316     for c in ordered:
317         if c in df_src.columns:
318             X[c] = df_src[c]
319         else:
320             X[c] = np.nan
321             msgs.append({"level": "warning", "code": "missing_train_col_created_nan",
322                         "col": c})
323
324     # Tipos exactos donde importa
325     if rf_meta.get("dtypes_expect", {}).get("Mes_lici") == "category" and "Mes_lici" in X.columns:
326         # Igual que en entrenamiento: Int->category (permitiendo NaN)
327         X["Mes_lici"] = pd.to_numeric(X["Mes_lici"], errors="coerce").astype("Int64").
328             astype("category")
329
330     # Numéricas a tipo numérico (coherente con imputer de mediana en el pipe)
331     for c in cols_num:
332         X[c] = pd.to_numeric(X[c], errors="coerce")

```

```

328
329     return X, msgs
330
331 def _ohe_unknowns_diagnostics(X: pd.DataFrame, rf_model) -> List[dict]:
332     """
333         Estima % de 'unknowns' por columna categórica frente al repertorio visto en train
334         (OneHotEncoder.categories_). Si hay unknowns, OHE pondrá todo a 0 para esa
335         columna.
336     """
337     msgs: List[dict] = []
338     try:
339         prep = rf_model.named_steps["prep"]
340         cat_cols: List[str] = []
341         oh = None
342         # Localizar el bloque categórico y su OneHotEncoder
343         for name, trans, cols in prep.transformers_:
344             if name == "cat":
345                 cat_cols = list(cols) if isinstance(cols, list) else []
346                 if hasattr(trans, "named_steps"):
347                     oh = trans.named_steps.get("oh") or trans.named_steps.get("ohe")
348                 else:
349                     oh = trans # Por si viene el OHE "plano" (raro)
350                     break
351         if oh is None or not hasattr(oh, "categories_"):
352             return msgs
353
354         # Alinear orden de columnas con categories_
355         n = min(len(cat_cols), len(oh.categories_))
356         for i in range(n):
357             col = cat_cols[i]
358             allowed = set(map(str, oh.categories_[i]))
359             s = X[col].astype(str)
360             unk = ~s.isin(allowed)
361             unk_rate = float(unk.mean()) if len(s) else 0.0
362             if unk_rate > 0.0:
363                 # muestra de desconocidos (máx 5)
364                 ex = s[unk].dropna().unique().tolist()[:5]
365                 msgs.append({
366                     "level": "warning", "code": "ohe_unknown_values",
367                     "col": col, "unknown_rate": round(unk_rate, 4),
368                     "examples": [str(x) for x in ex]
369                 })
370     except Exception as e:
371         msgs.append({"level": "info", "code": "ohe_diag_skipped", "reason": repr(e)})
372
373 def predict_rf_base(df_norm: pd.DataFrame, save_diag: bool = True, diag_tag: Optional[str] = None
374 ) -> Tuple[np.ndarray, np.ndarray, np.ndarray, List[str]]:
375     """
376         Adaptador ESTRICTO para RF base:
377         1) materializa alias EXACTOS usados en train (según rf_base_meta)
378         2) sub-conjunta columnas crudas en el ORDEN EXACTO del ColumnTransformer
379         3) castea tipos (Mes_lici -> category) y numéricas a float/int
380         4) diagnóstico de unknowns en OHE
381         Devuelve: (pred, conf, proba, labels) – igual interfaz que el adaptador previo.
382         Si save_diag=True, guarda un JSON con el diagnóstico y un CSV snapshot de X.
383     """
384     # 1) alias usados en train
385     df_used, m1 = _ensure_train_used_aliases(df_norm, rf_base_meta)
386     # 2) subset exacto + tipos
387     X, m2 = _subset_exact_order_for_pipe(df_used, rf_base_meta)
388     # 3) diagnóstico OHE unknowns
389     m3 = _ohe_unknowns_diagnostics(X, rf_base)
390
391     # 4) predicción
392     proba = rf_base.predict_proba(X)
393     labels = [str(c) for c in rf_base.named_steps["clf"].classes_]
394     pred_idx = np.argmax(proba, axis=1)
395     pred = np.array([labels[i] for i in pred_idx], dtype=object)
396     conf = np.array([proba[i, pred_idx[i]] for i in range(len(pred_idx))], dtype=float)
397

```

```

397
398     # 5) guardado diagnóstico
399     if save_diag:
400         ts = datetime.now().strftime("%Y%m%d %H%M%S")
401         tag = f"_ {diag_tag}" if diag_tag else ""
402         ddir = DIR_OUTPUTS / f"rf_base_diag_{ts}{tag}"
403         ddir.mkdir(parents=True, exist_ok=True)
404         # Snapshot del input exacto que entra al RF
405         X.to_csv(ddir / "rf_base_input_snapshot.csv", index=False)
406         # Meta/diagnóstico
407         diag_payload = {
408             "messages": m1 + m2 + m3,
409             "expected_raw_features": rf_base_meta.get("expected_raw_features", {}),
410             "alias_mapping_used_in_train": rf_base_meta.get(
411                 "alias_mapping_used_in_train", {}),
412             "notes": "unknown_rate>0 en alguna categórica indica valores no vistos en
413             train"
414         }
415         with open(ddir/"rf_base_diag.json", "w", encoding="utf-8") as f:
416             json.dump(diag_payload, f, indent=2, ensure_ascii=False)
417             print(f"RF base: diagnóstico guardado en {ddir}")
418
419     return pred, conf, proba, labels
420
421 """
422 =====
423 =
424 # RF_124_BASE+BC – Adaptador ESTRICTO (alias de train + orden + dtypes + duplicadas
425 # BC)
426 # Reemplaza o coloca tras tu Celda 4 (redefine predict_rf124)
427 # Requiere: rf_124, rf_124_meta, DIR_OUTPUTS ya cargados (Celda 0) y baja_comp
428 # (Series)
429 #
430 =====
431 =
432 from typing import Tuple, List, Dict, Optional
433 import numpy as np, pandas as pd, json
434 from datetime import datetime
435
436 def _ensure_train_used_aliases_124(df_norm: pd.DataFrame, rf124_meta: dict) -> Tuple[
437     pd.DataFrame, List[dict]]:
438     """
439         Materializa en df_norm las columnas EXACTAS con las que se entrenó el RF_124
440         (según
441         rf124_basebc_meta['alias_mapping_used_in_train']). Si faltan, intenta crearlas
442         desde
443         su 'preferred' o alias comunes (C_precio_p <-> Criterio_precio_p, Plazo_m <->
444         Plazo_m_c).
445     """
446     msgs: List[dict] = []
447     out = df_norm.copy()
448
449     alias_train: Dict[str, str] = rf124_meta.get("alias_mapping_used_in_train", {}) or
450     {}
451     fallback_alias = {
452         "C_precio_p": ["Criterio_precio_p"],
453         "Criterio_precio_p": ["C_precio_p"],
454         "Plazo_m_c": ["Plazo_m"],
455         "Plazo_m": ["Plazo_m_c"],
456     }
457
458     for preferred, used in alias_train.items():
459         if used in out.columns:
460             continue
461         if preferred in out.columns:
462             out[used] = out[preferred]
463             msgs.append({"level": "info", "code": "alias_materialized", "used": used, "from":
464             "preferred"})
465             continue

```

```

454     fb_cands = fallback_alias.get(used, []) + fallback_alias.get(preferred, [])
455     created = False
456     for alt in fb_cands:
457         if alt in out.columns:
458             out[used] = out[alt]
459             msgs.append({"level": "info", "code": "alias_fallback_materialized",
460                         "used": used, "from": alt})
461             created = True
462             break
463     if not created:
464         out[used] = np.nan
465         msgs.append({"level": "warning", "code": "alias_missing_filled_nan", "used": used})
466
467     return out, msgs
468
469 def _subset_exact_order_for_pipe_124(df_src: pd.DataFrame, rf124_meta: dict) -> Tuple[pd.DataFrame, List[dict]]:
470     """
471         Subconjunta en el ORDEN EXACTO las columnas crudas que el ColumnTransformer del
472         RF_124 espera.
473         Si falta alguna, la crea como NaN (y lo reporta). 'Mes_lici' se castea a category
474         si aplica.
475         Convierte numéricas a tipo numérico (el imputer de mediana está dentro del pipe).
476     """
477     msgs: List[dict] = []
478     exp: Dict[str, List[str]] = rf124_meta.get("expected_raw_features", {} or {})
479     cols_num = list(exp.get("num", []))
480     cols_cat = list(exp.get("cat", []))
481     ordered = cols_num + cols_cat
482
483     X = pd.DataFrame(index=df_src.index, columns=ordered, dtype="object")
484     for c in ordered:
485         if c in df_src.columns:
486             X[c] = df_src[c]
487         else:
488             X[c] = np.nan
489             msgs.append({"level": "warning", "code": "missing_train_col_created_nan",
490                         "col": c})
491
492     # Tipos exactos
493     if rf124_meta.get("alias_mapping_used_in_train", {}) is not None and "Mes_lici" in
494         X.columns:
495         X["Mes_lici"] = pd.to_numeric(X["Mes_lici"], errors="coerce").astype("Int64").
496             astype("category")
497     for c in cols_num:
498         X[c] = pd.to_numeric(X[c], errors="coerce")
499
500     return X, msgs
501
502 def _attach_baja_comp_dups(X: pd.DataFrame, baja_comp: pd.Series, rf124_meta: dict) ->
503     Tuple[pd.DataFrame, List[dict]]:
504     """
505         Inserta 'baja_comp' (categoría) y sus duplicadas EXACTAS tal como se usaron en
506         entrenamiento
507         (rf124_meta['baja_comp_dup_names']).
508     """
509     msgs: List[dict] = []
510     X = X.copy()
511     # base
512     X["baja_comp"] = pd.Categorical(baja_comp)
513     if "baja_comp" not in list(X.columns):
514         msgs.append({"level": "warning", "code": "baja_comp_missing_in_expected", "hint": "
515                         "Se añadió 'baja_comp' aunque no estuviera en expected_raw_features"})
516     # duplicadas
517     dup_names = rf124_meta.get("baja_comp_dup_names", [])
518     for dn in dup_names:
519         X[dn] = X["baja_comp"]
520     # Nota: si el pipe esperaba las duplicadas en expected_raw_features.cat, ya están
521     # presentes por nombre exacto.
522
523     return X, msgs
524
525 def _ohe_unknowns_diagnostics_124(X: pd.DataFrame, rf124_model) -> List[dict]:

```

```

514
515     """
516     Diagnóstico de unknowns por columna categórica frente al repertorio de OHE del
517     especialista.
518     """
519     msgs: List[dict] = []
520     try:
521         prep = rf124_model.named_steps["prep"]
522         cat_cols: List[str] = []
523         oh = None
524         for name, trans, cols in prep.transformers_:
525             if name == "cat":
526                 cat_cols = list(cols) if isinstance(cols, list) else []
527                 if hasattr(trans, "named_steps"):
528                     oh = trans.named_steps.get("oh") or trans.named_steps.get("ohe")
529                 else:
530                     oh = trans
531             break
532         if oh is None or not hasattr(oh, "categories_"):
533             return msgs
534
535         n = min(len(cat_cols), len(oh.categories_))
536         for i in range(n):
537             col = cat_cols[i]
538             allowed = set(map(str, oh.categories_[i]))
539             s = X[col].astype(str)
540             unk = ~s.isin(allowed)
541             rate = float(unk.mean()) if len(s) else 0.0
542             if rate > 0.0:
543                 ex = s[unk].dropna().unique().tolist()[:5]
544                 msgs.append({"level": "warning", "code": "ohe_unknown_values", "col": col,
545                             "unknown_rate": round(rate, 4), "examples": [str(x) for x in ex]})
546     except Exception as e:
547         msgs.append({"level": "info", "code": "ohe_diag_skipped", "reason": repr(e)})
548     return msgs
549
550 def predict_rf124(df_norm: pd.DataFrame, baja_comp: pd.Series, save_diag: bool = True,
551                    diag_tag: Optional[str] = None
552                               ) -> Tuple[np.ndarray, np.ndarray, np.ndarray, List[str]]:
553     """
554     Adaptador ESTRICTO para RF_124:
555     1) materializa alias EXACTOS usados en train (según rf_124_meta)
556     2) sub-conjunta columnas crudas en el ORDEN EXACTO del ColumnTransformer
557     3) castea tipos (Mes_lici -> category) y numéricas a float/int
558     4) inserta 'baja_comp' + duplicadas EXACTAS (nombres del meta)
559     5) diagnóstico de unknowns en OHE
560     Devuelve: (pred, conf, proba, labels) – misma interfaz.
561     """
562
563     # 1) alias de train
564     df_used, m1 = _ensure_train_used_aliases_124(df_norm, rf_124_meta)
565     # 2) subset exacto + tipos
566     X, m2 = _subset_exact_order_for_pipe_124(df_used, rf_124_meta)
567     # 3) baja_comp y duplicadas EXACTAS
568     X, m3 = _attach_baja_comp_dups(X, baja_comp, rf_124_meta)
569     # 4) diagnóstico OHE
570     m4 = _ohe_unknowns_diagnosticos_124(X, rf_124)
571
572     # 5) predicción
573     proba = rf_124.predict_proba(X)
574     labels = [str(c) for c in rf_124.named_steps["clf"].classes_]
575     pred_idx = np.argmax(proba, axis=1)
576     pred = np.array([labels[i] for i in pred_idx], dtype=object)
577     conf = np.array([proba[i, pred_idx[i]] for i in range(len(pred_idx))], dtype=float)
578
579     # 6) guardado diagnóstico
580     if save_diag:
581         ts = datetime.now().strftime("%Y%m%d %H%M%S")
582         tag = f"_{{diag_tag}}" if diag_tag else ""
583         ddir = DIR_OUTPUTS / f"rf124_basebc_diag_{ts}{tag}"
584         ddir.mkdir(parents=True, exist_ok=True)
585         X.to_csv(ddir / "rf124_input_snapshot.csv", index=False)
586         with open(ddir / "rf124_diag.json", "w", encoding="utf-8") as f:

```

```

582
583     json.dump({
584         "messages": m1 + m2 + m3 + m4,
585         "expected_raw_features": rf_124_meta.get("expected_raw_features", {}),
586         "alias_mapping_used_in_train": rf_124_meta.get(
587             "alias_mapping_used_in_train", {}),
588             "baja_comp_dup_names": rf_124_meta.get("baja_comp_dup_names", []),
589             "notes": "unknowns>0 en categóricas del especialista degradan su
590             margen vs RF base y reducen flips."
591     }, f, indent=2, ensure_ascii=False)
592     print(f"RF_124: diagnóstico guardado en {ddir}")
593
594     return pred, conf, proba, labels
595
596 """Celda 5 - Ensemble "seguro" (scope {1,2,4}, margen desde ensemble_meta.json) +
597 segundo cluster"""
598
599 # =====
600 # ENSEMBLE - BASE vs 124 (scope {1,2,4}) + segundo cluster
601 # =====
602
603 def apply_ensemble_from_parts(pred_base, conf_base, proba_base, labels_base,
604                                pred_124, conf_124, proba_124, labels_124,
605                                margin_min=None):
606
607     if margin_min is None:
608         margin_min = float(ensemble_meta.get("margin_min", 0.01))
609
610     idx_b = {lb:i for i,lb in enumerate(labels_base)}
611     idx_s = {lb:i for i,lb in enumerate(labels_124)}
612
613     final_pred = pred_base.copy()
614     fuente = np.array(["rf_base"]*len(final_pred), dtype=object)
615     conf_final = conf_base.copy()
616     flips = np.zeros(len(final_pred), dtype=bool)
617
618     # Flip SOLO si la clase base está en {1,2,4} y el especialista gana por margen
619     for i in range(len(final_pred)):
620         base_cls = str(pred_base[i])
621         alt_cls = str(pred_124[i])
622         if base_cls not in CLASSES_124:
623             continue
624         if (alt_cls in CLASSES_124) and (alt_cls != base_cls):
625             mb = proba_base[i, idx_b[base_cls]]
626             ma = proba_124[i, idx_s[alt_cls]]
627             if (ma - mb) >= margin_min:
628                 final_pred[i] = alt_cls
629                 fuente[i] = "rf_124_basebc"
630                 conf_final[i] = ma
631                 flips[i] = True
632
633     # Segundo cluster (si conf_final < 0.70) usando el modelo decisor
634     segundo = np.array([None]*len(final_pred), dtype=object)
635     conf2 = np.array([np.nan]*len(final_pred), dtype=float)
636     for i in range(len(final_pred)):
637         if fuente[i] == "rf_124_basebc":
638             probs, labs = proba_124[i], labels_124
639         else:
640             probs, labs = proba_base[i], labels_base
641             order = np.argsort(probs)[::-1]
642             top1 = labs[order[0]]
643             p1 = probs[order[0]]
644             # coherencia
645             if str(top1) != str(final_pred[i]):
646                 top1 = final_pred[i]
647                 p1 = conf_final[i]
648                 order = [np.where(np.array(labs)==top1)[0][0]] + [j for j in order if labs[j]!=top1]
649             if p1 < 0.70 and len(order) > 1:
650                 sc = labs[order[1]]; p2 = probs[order[1]]
651                 segundo[i], conf2[i] = sc, float(p2)
652
653     out = pd.DataFrame({
654         "pred": final_pred,
655         "confianza_pred": conf_final,

```

```

650         "fuente_confianza": fuente,
651         "segundo_cluster": segundo,
652         "confianza_segundo": conf2,
653         "flip_ensemble": flips
654     })
655     return out
656
657 """Celda 6 - kNN similares (NUM-ONLY con su pipeline exacta) + catálogo"""
658
659 # =====
660 # kNN - similares (NUM-ONLY) con pipeline exacta + catálogo
661 # =====
662 def knn_similares(df_norm: pd.DataFrame, topk=3, sim_thr=60.0):
663     pre_num_features = knn_params["pre_num_features"]
664     p95_nn1_train = float(knn_params["p95_nn1_train"])
665     # Construir matriz NUM-ONLY en el ORDEN exacto y transformar con knn_num_pipeline
666     X_num = pd.DataFrame(index=df_norm.index)
667     for c in pre_num_features:
668         if c in df_norm.columns:
669             X_num[c] = pd.to_numeric(df_norm[c], errors="coerce")
670         elif c == "C_precio_p" and "Criterio_precio_p" in df_norm.columns:
671             X_num[c] = pd.to_numeric(df_norm["Criterio_precio_p"], errors="coerce")
672         else:
673             X_num[c] = np.nan
674     Zq = knn_num_pipe.transform(X_num)
675
676     dist, idx = nn_all.kneighbors(Zq, n_neighbors=max(topk, 1), return_distance=True)
677     def sim_percent(d):
678         return max(0.0, 1.0 - float(d)/p95_nn1_train)*100.0 if p95_nn1_train>0 else
679         0.0
680
681     out_rows = []
682     for i in range(len(df_norm)):
683         pairs = [(float(dist[i,j]), int(idx[i,j])) for j in range(dist.shape[1])]
684         sims = [(j, sim_percent(d)) for (d,j) in pairs]
685         keep = [(j, s) for (j,s) in sims if s >= sim_thr]
686         warn = False
687         if len(keep) == 0:
688             keep = [max(sims, key=lambda t:t[1])]
689             warn = True
690         keep = sorted(keep, key=lambda t:t[1], reverse=True)[:topk]
691         for rank, (j,s) in enumerate(keep, start=1):
692             row = {"input_row": i, "rank": rank, "knn_row": j, "sim%": round(float(s),
693             2)}
694             if catalog is not None:
695                 meta = catalog.iloc[j]
696                 row["Cluster_6"] = str(meta.get("Cluster_6", ""))
697                 if id_col: row["id"] = meta.get(id_col, f"row_{j}")
698                 for key in ["Presupuesto_licitacion_lote_c", "Plazo_m",
699                             "Criterio_precio_p", "N_ofertantes"]:
700                     if key in catalog.columns:
701                         row[key] = meta.get(key, np.nan)
702                 if warn: row["warning"] = "No hay similares con sim>=60%; se muestra el
703                 más cercano."
704             out_rows.append(row)
705     return pd.DataFrame(out_rows)
706
707 """Celda 7 - POST stats para pred (y opcional segundo_cluster)"""
708
709 # =====
710 # POST stats - medias y std por cluster (pred y segundo)
711 # =====
712 def post_stats_for_clusters(pred_df: pd.DataFrame):
713     if catalog is None or "Cluster_6" not in catalog.columns:
714         return pd.DataFrame(), pd.DataFrame()
715     groups = catalog.groupby(catalog["Cluster_6"].astype(str))
716     def _stats_for(cluster):
717         g = groups.get_group(str(cluster)) if str(cluster) in groups.indices else None
718         rec = {}
719         for pref, col in POST_MAP.items():
720             if g is None or col is None:
721                 rec[f"{pref}_mean"] = np.nan; rec[f"{pref}_std"] = np.nan

```

```

718
719     else:
720         rec[f"{pref}_mean"] = float(g[col].mean(numeric_only=True))
721         rec[f"{pref}_std"] = float(g[col].std(numeric_only=True))
722     return rec
723
724 post1_rows, post2_rows = [], []
725 for i, row in pred_df.reset_index(drop=True).iterrows():
726     cl = row["pred"]
727     rec1 = {"input_row": i, "cluster": str(cl)}
728     rec1.update(_stats_for(cl))
729     post1_rows.append(rec1)
730     sc = row.get("segundo_cluster", None)
731     if pd.notna(sc):
732         rec2 = {"input_row": i, "cluster": str(sc)}
733         rec2.update(_stats_for(sc))
734         post2_rows.append(rec2)
735 return pd.DataFrame(post1_rows), pd.DataFrame(post2_rows)

736 """Celda 8 - (Opcional) Coherencia geométrica via centroides (separada)"""
737
738 # =====
739 # (OPCIONAL) Coherencia geométrica por centroides
740 # =====
741 def coherencia_geometrica(df_norm: pd.DataFrame):
742     Z = preprocess.transform(df_norm) # espacio PRE (imputación+OHE+escala)
743     classes = [str(c) for c in centroides["classes"]]
744     cents = {str(k): np.array(v, dtype=float) for k,v in centroides["centroids_pre"].items()}
745     p95_intra= {str(k): float(v) for k,v in centroides.get("p95_intra", {}).items()}
746
747     out=[]
748     for i in range(Z.shape[0]):
749         z = Z[i]
750         d = {c: float(np.linalg.norm(z - cents[c])) for c in classes}
751         items = sorted(d.items(), key=lambda kv: kv[1])
752         cmin, dmin = items[0]
753         d2 = items[1][1] if len(items)>1 else np.nan
754         margin = float(d2-dmin) if np.isfinite(d2) else np.nan
755         p95 = p95_intra.get(cmin, np.nan)
756         coher = (1.0 - dmin/p95)*100.0 if (p95 and p95>0) else np.nan
757         out.append({"min_cluster_geom": cmin, "min_dist": dmin, "geom_margin": margin,
758                     "coherencia_geom": coher})
759
760 return pd.DataFrame(out)

761 """Celda 9 - Orquestación: de inputs → todas las salidas (CSV por lotes)"""
762
763 # =====
764 # ORQUESTADOR - batch: inputs CSV → outputs CSV
765 # =====
766 def run_batch(csv_path: Path, topk=3, sim_thr=60.0, add_geom=False):
767     df_in = pd.read_csv(csv_path) if csv_path.suffix.lower() == ".csv" else pd.read_excel(csv_path)
768     df_norm, _ = validate_normalize_df(df_in, strict=False)
769
770     # 1) baja_comp (XGB)
771     baja, p_bajo = predict_baja_comp_xgb(df_norm)
772
773     # 2) RF base
774     pred_b, conf_b, proba_b, labs_b = predict_rf_base(df_norm)
775
776     # 3) RF_124
777     pred_s, conf_s, proba_s, labs_s = predict_rf124(df_norm, baja)
778
779     # 4) Ensemble
780     ens = apply_ensemble_from_parts(pred_b, conf_b, proba_b, labs_b,
781                                     pred_s, conf_s, proba_s, labs_s,
782                                     margin_min=None)
783     # Enriquecer con baja_comp
784     ens.insert(0, "baja_comp", baja.values)
785     ens.insert(1, "p_bajo", p_bajo.values)
786
787     # 5) Similares kNN

```

```

787     sims = knn_similares(df_norm, topk=topk, sim_thr=sim_thr)
788
789     # 6) POST stats
790     post1, post2 = post_stats_for_clusters(ens)
791
792     # 7) (opcional) coherencia geométrica
793     geom = None
794     if add_geom:
795         geom = coherencia_geometrica(df_norm)
796
797     # Guardados
798     ts = pd.Timestamp.now().strftime("%Y%m%d_%H%M%S")
799     out_dir = DIR_OUTPUTS / f"pred_{ts}"
800     out_dir.mkdir(parents=True, exist_ok=True)
801
802     # Adjuntar ID si está
803     id_cand = next((c for c in ["Identificador", "identificador", "ID", "Id", "id_lote",
804                      "id_contrato"] if c in df_in.columns), None)
805     main = ens.copy()
806     if id_cand:
807         main.insert(0, "Identificador", df_in[id_cand].values[:len(main)])
808
809     # Añadir POST stats (pred y segundo) en columnas
810     if not post1.empty:
811         sel = []
812         for i in range(len(main)):
813             r = post1[post1["input_row"]==i]
814             rec={"input_row":i}
815             if len(r):
816                 for pref in POST_PREFS:
817                     rec[f"POST_{pref}_mean"] = float(r.iloc[0].get(f"{pref}_mean", np.nan))
818                     rec[f"POST_{pref}_std"] = float(r.iloc[0].get(f"{pref}_std", np.nan))
819             sel.append(rec)
820         main = main.join(pd.DataFrame(sel).set_index("input_row"), how="left")
821     if not post2.empty:
822         sel2=[]
823         for i in range(len(main)):
824             r = post2[post2["input_row"]==i]
825             rec={"input_row":i}
826             if len(r):
827                 for pref in POST_PREFS:
828                     rec[f"POST2_{pref}_mean"] = float(r.iloc[0].get(f"{pref}_mean", np.nan))
829                     rec[f"POST2_{pref}_std"] = float(r.iloc[0].get(f"{pref}_std", np.nan))
830             sel2.append(rec)
831         main = main.join(pd.DataFrame(sel2).set_index("input_row"), how="left")
832
833     main.to_csv(out_dir/"predicciones.csv", index=False)
834     sims.to_csv(out_dir/"similares_topK.csv", index=False)
835     post1.to_csv(out_dir/"post_stats_pred.csv", index=False)
836     if not post2.empty:
837         post2.to_csv(out_dir/"post_stats_segundo.csv", index=False)
838     if geom is not None:
839         geom.to_csv(out_dir/"coherencia_geom.csv", index=False)
840
841     print(f"✓ Batch guardado en: {out_dir}")
842     return {"main": main, "sims": sims, "post1": post1, "post2": post2, "geom": geom}
843
844     """Celda - predict_one para GPT (JSON)"""
845
846     # =====
847     # API para GPT - predict_one(input_dict, k=3, sim_thr=60)
848     # =====
849     def predict_one(input_dict: dict, topk: int = 3, sim_thr: float = 60.0,
850                      include_geom: bool = False, include_validation_messages: bool = True):
851         """
852             Entrada: dict con SOLO variables PRE (se aceptan alias C_precio_p/Plazo_m_c,
853             Mes_lici texto, etc.)
854             Salida: dict JSON con:

```

```

853     - baja_comp, p_bajo
854     - pred, confianza_pred, fuente_confianza, (si <70%) segundo_cluster,
855       confianza_segundo, flip_ensemble
856     - post_stats_pred {var: {mean, std}} y opcional post_stats_segundo
857     - similares (hasta 3 con sim>=60%; si no, 1 con warning)
858     - (opcional) coherencia_geom {min_cluster_geom, min_dist, geom_margin,
859       coherencia_geom}
860     - (opcional) validation_messages (para revisar normalizaciones/fallbacks)
861 """
862 # 0) Validación / normalización (común)
863 df_raw = pd.DataFrame([input_dict])
864 df_norm, messages = validate_normalize_df(df_raw, strict=False)
865
866 # 1) baja_comp (XGB)
867 baja, p_bajo = predict_baja_comp_xgb(df_norm)
868
869 # 2) RF base
870 pred_b, conf_b, proba_b, labs_b = predict_rf_base(df_norm)
871
872 # 3) RF_124 (especialista)
873 pred_s, conf_s, proba_s, labs_s = predict_rf124(df_norm, baja)
874
875 # 4) Ensemble (scope {1,2,4} + margen)
876 ens = apply_ensemble_from_parts(pred_b, conf_b, proba_b, labs_b,
877                                 pred_s, conf_s, proba_s, labs_s,
878                                 margin_min=None)
879
880 # 5) kNN similares
881 sims_df = knn_similares(df_norm, topk=topk, sim_thr=sim_thr)
882 sims0 = sims_df[sims_df["input_row"]==0].sort_values("rank").to_dict(orient=
883 "records")
884
885 # 6) POST stats (pred y, si aplica, segundo)
886 post_pred_df, post_second_df = post_stats_for_clusters(ens)
887 post_pred = {}
888 if not post_pred_df.empty:
889     r = post_pred_df[post_pred_df["input_row"]==0]
890     if len(r):
891         post_pred = {pref: {"mean": float(r.iloc[0].get(f"{pref}_mean", np.nan)),
892                            "std": float(r.iloc[0].get(f"{pref}_std", np.nan))}}
893         for pref in POST_PREFS}
894 post_second = {}
895 if not post_second_df.empty and pd.notna(ens.loc[0, "segundo_cluster"]):
896     r2 = post_second_df[post_second_df["input_row"]==0]
897     if len(r2):
898         post_second = {pref: {"mean": float(r2.iloc[0].get(f"{pref}_mean", np.nan)),
899                           "std": float(r2.iloc[0].get(f"{pref}_std", np.nan))}}
900         for pref in POST_PREFS}
901
902 # 7) (Opcional) coherencia geométrica
903 geom = None
904 if include_geom:
905     geom_df = coherencia_geometrica(df_norm)
906     if not geom_df.empty:
907         g = geom_df.iloc[0].to_dict()
908         # redondeos amigables
909         g["min_dist"] = float(g.get("min_dist", np.nan))
910         g["geom_margin"] = float(g.get("geom_margin", np.nan)) if pd.notna(g.get(
911             "geom_margin", np.nan)) else None
912         g["coherencia_geom"] = float(g.get("coherencia_geom", np.nan)) if pd.notna(
913             g.get("coherencia_geom", np.nan)) else None
914         geom = g
915
916 # 8) Ensamblar respuesta JSON
917 rec = {
918     "baja_comp": int(ens.loc[0, "baja_comp"]),
919     "p_bajo": float(ens.loc[0, "p_bajo"]),
920     "pred": str(ens.loc[0, "pred"]),
921     "confianza_pred": float(ens.loc[0, "confianza_pred"]),
922     "fuente_confianza": str(ens.loc[0, "fuente_confianza"]),
923     "flip_ensemble": bool(ens.loc[0, "flip_ensemble"]),
924     "segundo_cluster": str(ens.loc[0, "segundo_cluster"]),
925     "coherencia_geom": float(ens.loc[0, "coherencia_geom"]),
926     "min_dist": float(ens.loc[0, "min_dist"]),
927     "geom_margin": float(ens.loc[0, "geom_margin"]),
928     "similarity": float(ens.loc[0, "similarity"])
929 }

```

```

918     "similares": sims0,
919     "post_stats_pred": post_pred
920   }
921   # segundo cluster si aplica
922   if pd.notna(ens.loc[0, "segundo_cluster"]):
923     rec["segundo_cluster"] = str(ens.loc[0, "segundo_cluster"])
924     rec["confianza_segundo"] = float(ens.loc[0, "confianza_segundo"])
925     if post_second:
926       rec["post_stats_segundo"] = post_second
927   # geom opcional
928   if include_geom and geom is not None:
929     rec["coherencia_geom"] = geom
930   # mensajes de validación opcionales
931   if include_validation_messages:
932     # Filtrar mensajes relevantes a la fila 0 y globales ("*")
933     msgs0 = [m for m in messages if m.get("row") in (0, "*")]
934     rec["validation_messages"] = msgs0
935
936   # Info de ensemble (from→to) si hubo flip
937   if rec["flip_ensemble"]:
938     rec["ensemble_info"] = {
939       "from": str(pred_b[0]),
940       "to": str(pred_s[0]),
941       "margin_delta": float(
942         (proba_s[0, list(labs_s).index(str(pred_s[0]))]) if str(pred_s[0]) in
943         labs_s else np.nan)
944       - (proba_b[0, list(labs_b).index(str(pred_b[0]))]) if str(pred_b[0]) in
945         labs_b else np.nan)
946     )
947   }
948
949   return rec
950
951 """
952 =====
953 # PRUEBA GLOBAL - 20 casos TEST + 2 CSV + evaluación completa (TEST entero)
954 #
955 =====
956
957 from sklearn.metrics import confusion_matrix, accuracy_score
958
959 # Verificaciones mínimas
960 assert callable(validate_normalize_df), "Ejecuta primero las celdas del pipeline de
961 producción."
962 assert callable(predict_baja_comp_xgb) and callable(predict_rf_base) and callable(
963 predict_rf124)
964 assert callable(apply_ensemble_from_parts) and callable(knn_similares) and callable(
965 post_stats_for_clusters)
966
967 # 0) Cargar dataset de inputs (TEST)
968 cand_paths = [
969   DIR_INPUTS/"df_train_test_v8.xlsx",
970   DIR_INPUTS/"df_test_train_v8.xlsx",
971   DIR_INPUTS/"df_train_test_v8.csv",
972   DIR_INPUTS/"df_test_train_v8.csv",
973 ]
974 data_path = next((p for p in cand_paths if p.exists()), None)
975 assert data_path is not None, f"X No se encontró dataset en: {cand_paths}"
976 df_all = pd.read_excel(data_path) if data_path.suffix.lower() in [".xlsx", ".xls"] else
977   pd.read_csv(data_path)
978 assert "Dataset" in df_all.columns, "X Falta columna 'Dataset' (TRAIN/TEST)"
979 df_all["Dataset"] = df_all["Dataset"].astype(str).str.upper()
980 df_test = df_all[df_all["Dataset"]=="TEST"].copy().reset_index(drop=True)
981 assert len(df_test)>0, "X El split TEST está vacío."
982
983 # 1) Selección reproducible de 20 casos
984 rng = np.random.RandomState(42)
985 idx20 = rng.choice(len(df_test), size=min(20, len(df_test)), replace=False)
986 df20 = df_test.iloc[idx20].reset_index(drop=True)

```

```

980
981 # 2) Predicción para los 20 casos
982 df_norm_20, _ = validate_normalize_df(df20, strict=False)
983 baja20, p_bajo20 = predict_baja_comp_xgb(df_norm_20)
984 pb_pred20, pb_conf20, pb_proba20, pb_labs = predict_rf_base(df_norm_20)
985 ps_pred20, ps_conf20, ps_proba20, ps_labs = predict_rf124(df_norm_20, baja20)
986 ens20 = apply_ensemble_from_parts(pb_pred20, pb_conf20, pb_proba20, pb_labs,
987                                     ps_pred20, ps_conf20, ps_proba20, ps_labs,
988                                     margin_min=None)
989 ens20.insert(0, "baja_comp", baja20.values)
990 ens20.insert(1, "p_bajo", p_bajo20.values)
991
992 # POST stats pred/segundo
993 post1_20, post2_20 = post_stats_for_clusters(ens20)
994
995 # Preparar salida de 20 (outputs genéricos de la app)
996 id_col = next((c for c in ["Identificador", "identificador", "ID", "Id", "id_lote",
997 "id_contrato"] if c in df20.columns), None)
998 main20 = ens20.copy()
999 if id_col:
1000     main20.insert(0, "Identificador", df20[id_col].values)
1001
1002 # Anexar POST stats (pred y segundo) en columnas
1003 if not post1_20.empty:
1004     sel = []
1005     for i in range(len(main20)):
1006         r = post1_20[post1_20["input_row"]==i]
1007         rec={"input_row":i}
1008         if len(r):
1009             for pref in POST_PREFS:
1010                 rec[f"POST_{pref}_mean"] = float(r.iloc[0].get(f"{pref}_mean", np.nan))
1011                 rec[f"POST_{pref}_std"] = float(r.iloc[0].get(f"{pref}_std", np.nan))
1012         sel.append(rec)
1013 main20 = main20.join(pd.DataFrame(sel).set_index("input_row"), how="left")
1014
1015 if not post2_20.empty:
1016     sel2 = []
1017     for i in range(len(main20)):
1018         r = post2_20[post2_20["input_row"]==i]
1019         rec={"input_row":i}
1020         if len(r):
1021             for pref in POST_PREFS:
1022                 rec[f"POST2_{pref}_mean"] = float(r.iloc[0].get(f"{pref}_mean", np.nan))
1023                 rec[f"POST2_{pref}_std"] = float(r.iloc[0].get(f"{pref}_std", np.nan))
1024         sel2.append(rec)
1025 main20 = main20.join(pd.DataFrame(sel2).set_index("input_row"), how="left")
1026
1027 # 3) Construir CSV de comparación POST pred vs real (20 casos)
1028 # Resolver mapping de columnas reales en df_all
1029 def resolve_post_columns_df(df_cols):
1030     mapping={}
1031     for pref, cands in POST_PREFS.items():
1032         found = next((c for c in cands if c in df_cols), None)
1033         mapping[pref] = found
1034     return mapping
1035 POST_MAP_DF = resolve_post_columns_df(set(df_all.columns))
1036
1037 rows = []
1038 for i in range(len(df20)):
1039     rec = {}
1040     if id_col: rec["Identificador"] = df20.loc[i, id_col]
1041     rec["pred_cluster"] = str(main20.loc[i, "pred"])
1042     rec["confianza_pred"] = float(main20.loc[i, "confianza_pred"])
1043     rec["baja_comp_pred"] = int(main20.loc[i, "baja_comp"])
1044     rec["p_bajo"] = float(main20.loc[i, "p_bajo"])
1045     # POST pred (medias del cluster)
1046     for pref in POST_PREFS:
1047         rec[f"{pref}_pred_mean"] = float(main20.loc[i, f"POST_{pref}_mean"]) if

```

```

1047     f"POST_{pref}_mean" in main20.columns else np.nan
1048     rec[f"{pref}_pred_std"] = float(main20.loc[i, f"POST_{pref}_std"]) if
1049     f"POST_{pref}_std" in main20.columns else np.nan
1050
1051 # POST reales
1052 for pref, col in POST_MAP_DF.items():
1053     rec[f"{pref}_real"] = (None if col is None else df20.loc[i, col])
1054
1055 # baja_comp_real derivada (si hay N_ofertantes real)
1056 if POST_MAP_DF["N_ofertantes"]:
1057     nreal = df20.loc[i, POST_MAP_DF["N_ofertantes"]]
1058     rec["baja_comp_real"] = (int(nreal <= 4) if pd.notna(nreal) else np.nan)
1059 else:
1060     rec["baja_comp_real"] = np.nan
1061
1062 # cluster real
1063 rec["cluster_real"] = (str(df20.loc[i, "Cluster_6"]) if "Cluster_6" in df20.
1064 columns else None)
1065 rows.append(rec)
1066 targets20 = pd.DataFrame(rows)
1067
1068 # 4) Evaluación del ENSEMBLE en TODO el TEST
1069 df_norm_test, _ = validate_normalize_df(df_test, strict=False)
1070 baja_t, p_bajo_t = predict_baja_comp_xgb(df_norm_test)
1071 pb_pred_t, pb_conf_t, pb_proba_t, pb_labs_t = predict_rf_base(df_norm_test)
1072 ps_pred_t, ps_conf_t, ps_proba_t, ps_labs_t = predict_rf124(df_norm_test, baja_t)
1073 ens_t = apply_ensemble_from_parts(pb_pred_t, pb_conf_t, pb_proba_t, pb_labs_t,
1074                                     ps_pred_t, ps_conf_t, ps_proba_t, ps_labs_t,
1075                                     margin_min=None)
1076
1077 y_true = df_test["Cluster_6"].astype(str).values
1078 y_pred = ens_t["pred"].astype(str).values
1079 labels_sorted = sorted(pd.unique(pd.concat([pd.Series(y_true), pd.Series(y_pred)])), key=lambda s:int(s))
1080 cm = confusion_matrix(y_true, y_pred, labels=labels_sorted)
1081 acc_global = accuracy_score(y_true, y_pred)
1082 acc_by_cls = []
1083 for cl in labels_sorted:
1084     m = (y_true == cl)
1085     acc_c = float((y_pred[m]==y_true[m]).mean()) if m.any() else np.nan
1086     acc_by_cls.append({"cluster": cl, "accuracy": acc_c})
1087 acc_by_cls = pd.DataFrame(acc_by_cls)
1088
1089 # 5) Guardados
1090 ts = pd.Timestamp.now().strftime("%Y%m%d_%H%M%S")
1091 out_dir = DIR_OUTPUTS / f"test_global_{ts}"
1092 out_dir.mkdir(parents=True, exist_ok=True)
1093
1094 main20.to_csv(out_dir/"predicciones_20.csv", index=False)
1095 targets20.to_csv(out_dir/"targets_pred_vs_real_20.csv", index=False)
1096 pd.DataFrame(cm, index=[f"real_{c}" for c in labels_sorted], columns=[f"pred_{c}" for
1097 c in labels_sorted]).to_csv(out_dir/"cm_test.csv")
1098 acc_by_cls.to_csv(out_dir/"accuracy_por_cluster.csv", index=False)
1099 with open(out_dir/"summary_eval_test.json", "w", encoding="utf-8") as f:
1100     json.dump({"accuracy_global": float(acc_global), "labels": labels_sorted}, f,
1101     indent=2, ensure_ascii=False)
1102
1103 print("✓ Prueba global completada")
1104 print(f"→ {out_dir/'predicciones_20.csv'}")
1105 print(f"→ {out_dir/'targets_pred_vs_real_20.csv'}")
1106 print(f"→ {out_dir/'cm_test.csv'} | {out_dir/'accuracy_por_cluster.csv'} |
1107 acc_global={acc_global:.4f}")

```