

```

1  ##### SCRIPT 18 - PYTHON
2  # PIPELINE DE EXPORTACION DE LA TRANSFORMACION DE LOS MODELOS
3  # A ARCHIVOS TEXTO-MATRIZ
4  #
5  # =====
6  # -*- coding: utf-8 -*-
7  #
8  # Celda 0 - Setup + rutas + comprobación de entorno
9  #
10 import os, json, math, datetime, textwrap
11 from pathlib import Path
12
13 import numpy as np
14 import pandas as pd
15 import joblib
16
17 # Opcional: XGBoost, si está disponible en el entorno
18 try:
19     import xgboost as xgb
20     HAS_XGB = True
21 except Exception:
22     HAS_XGB = False
23
24 # Rutas (ajusta si cambias la estructura en Drive)
25 MODEL_DIR = Path("/content/drive/MyDrive/_Pipeline_produccion_prediccion/modelos")
26 TEST_XLSX = Path(
27     "/content/drive/MyDrive/_Pipeline_produccion_prediccion/inputs/df_test_train_v8.xlsx")
28 OUT_DIR = Path(
29     "/content/drive/MyDrive/_Pipeline_produccion_prediccion/modelos_portable")
30
31 OUT_DIR.mkdir(parents=True, exist_ok=True)
32
33 print("MODEL_DIR:", MODEL_DIR)
34 print("TEST_XLSX:", TEST_XLSX)
35 print("OUT_DIR : ", OUT_DIR)
36
37 assert MODEL_DIR.exists(), f"✗ MODEL_DIR no existe: {MODEL_DIR}"
38 assert TEST_XLSX.exists(), f"✗ No se encuentra el Excel de test: {TEST_XLSX}"
39
40 print("✓ Rutas verificadas.")
41 print("✓ XGBoost disponible: " if HAS_XGB else "⚠️ XGBoost NO disponible; limitar export.")
42
43 """Celda 1 - Carga de artefactos originales"""
44
45 # Celda 1 - Carga de artefactos originales
46
47 paths = {
48     "preprocess_meta": MODEL_DIR/"preprocess_meta.json",
49     "preprocess": MODEL_DIR/"preprocess_pipeline.pkl",
50     "centroides": MODEL_DIR/"centrodes_pre.json",
51
52     "rf_base": MODEL_DIR/"rf_base.pkl",
53     "rf_base_meta": MODEL_DIR/"rf_base_meta.json",
54
55     "rf_124": MODEL_DIR/"rf_124_basebc.pkl",
56     "rf_124_meta": MODEL_DIR/"rf_124_basebc_meta.json",
57
58     "xgb_pack": MODEL_DIR/"model_Noferta_XGB_v8.pkl",
59     "xgb_meta": MODEL_DIR/"xgb_meta.json",
60
61     "knn_train": MODEL_DIR/"knn_pre_train.joblib",
62     "knn_all": MODEL_DIR/"knn_pre_all.joblib",
63     "knn_params": MODEL_DIR/"knn_params.json",
64     "knn_num_pipe": MODEL_DIR/"knn_num_pipeline.joblib",
65     "knn_catalog": MODEL_DIR/"knn_catalog_all.parquet",
66     "knn_catalog_meta": MODEL_DIR/"knn_catalog_meta.json",
67
68     "ensemble_meta": MODEL_DIR/"ensemble_meta.json",
69 }
70
71 # Checks mínimos (catálogo opcional pero recomendable)

```

```

70  for k, p in paths.items():
71      if k in ["knn_catalog", "knn_catalog_meta"]:
72          continue
73      assert p.exists(), f"✗ Falta {k}: {p}"
74
75  # Cargas
76 preprocess_meta = json.loads(paths["preprocess_meta"].read_text(encoding="utf-8"))
77 preprocess = joblib.load(paths["preprocess"])
78 centroides = json.loads(paths["centroides"].read_text(encoding="utf-8"))
79
80 rf_base = joblib.load(paths["rf_base"])
81 rf_base_meta = json.loads(paths["rf_base_meta"].read_text(encoding="utf-8"))
82
83 rf_124 = joblib.load(paths["rf_124"])
84 rf_124_meta = json.loads(paths["rf_124_meta"].read_text(encoding="utf-8"))
85
86 xgb_pack = joblib.load(paths["xgb_pack"])
87 xgb_meta = json.loads(paths["xgb_meta"].read_text(encoding="utf-8"))
88
89 nn_train = joblib.load(paths["knn_train"])
90 nn_all = joblib.load(paths["knn_all"])
91 knn_params = json.loads(paths["knn_params"].read_text(encoding="utf-8"))
92 knn_num_pipe = joblib.load(paths["knn_num_pipe"])
93
94 ensemble_meta = json.loads(paths["ensemble_meta"].read_text(encoding="utf-8"))
95
96 catalog, knn_catalog_meta, knn_id_col, knn_index_key = None, None, None, None
97 if paths["knn_catalog"].exists():
98     catalog = pd.read_parquet(paths["knn_catalog"])
99 if paths["knn_catalog_meta"].exists():
100    knn_catalog_meta = json.loads(paths["knn_catalog_meta"].read_text(encoding="utf-8"))
101    knn_id_col = knn_catalog_meta.get("id_col")
102    knn_index_key = knn_catalog_meta.get("index_key")
103
104 print("✓ Artefactos originales cargados correctamente.")
105 print("  RF base meta keys:", list(rf_base_meta.keys()))
106 print("  RF 124 meta keys : ", list(rf_124_meta.keys()))
107 print("  XGB meta keys   : ", list(xgb_meta.keys()))
108
109 """Celda 2 - Utilidades de serialización segura a JSON"""
110
111 # Celda 2 - Utilidades de serialización segura a JSON
112
113 def _to_serializable(obj):
114     """Convierte numpy/pandas a tipos básicos JSON (list, float, int, dict...)."""
115     if isinstance(obj, (np.floating,)):
116         return float(obj)
117     if isinstance(obj, (np.integer,)):
118         return int(obj)
119     if isinstance(obj, (np.bool_,)):
120         return bool(obj)
121     if isinstance(obj, (np.ndarray,)):
122         return obj.tolist()
123     if isinstance(obj, (pd.Series, pd.Index)):
124         return obj.tolist()
125     if isinstance(obj, pd.DataFrame):
126         return obj.to_dict(orient="list")
127     if isinstance(obj, (dict, list, tuple)):
128         if isinstance(obj, dict):
129             return {k: _to_serializable(v) for k, v in obj.items()}
130         else:
131             return [_to_serializable(v) for v in obj]
132     return obj
133
134 def save_json(obj, path: Path):
135     path = Path(path)
136     obj2 = _to_serializable(obj)
137     path.write_text(json.dumps(obj2, indent=2, ensure_ascii=False), encoding="utf-8")
138     print(f"Guardado JSON: {path}")
139     return path
140

```

```

141 """Celda 3 - Inspección del preprocesador PRE (coherencia geométrica)"""
142
143 # Celda 3 - Inspección del preprocesador PRE (coherencia geométrica)
144
145 from sklearn.compose import ColumnTransformer
146 from sklearn.pipeline import Pipeline
147 from sklearn.impute import SimpleImputer
148 from sklearn.preprocessing import OneHotEncoder, StandardScaler
149
150 assert isinstance(preprocess, ColumnTransformer), "✗ preprocess no es ColumnTransformer."
151
152 print("Transformers en preprocess:")
153 for name, trans, cols in preprocess.transformers_:
154     print(f" - {name}: {type(trans).__name__}, cols={len(cols)}")
155
156 def _get_ct_block(ct: ColumnTransformer, block_name: str):
157     """
158         Devuelve (transformer, lista_de_columnas) para el bloque block_name
159         dentro de un ColumnTransformer.
160     """
161     for name, trans, cols in ct.transformers_:
162         if name == block_name:
163             return trans, list(cols)
164     raise KeyError(f"Blöque '{block_name}' no encontrado en ColumnTransformer.")
165
166 def _find_step(pipe: Pipeline, candidate_names, expected_type):
167     """
168         Busca un paso en un Pipeline:
169         1) Primero por nombre (en orden de candidate_names).
170         2) Si no, por tipo (expected_type).
171
172         Lanza AssertionError si no encuentra nada coherente.
173     """
174     # 1) Por nombre
175     for nm in candidate_names:
176         if nm in pipe.named_steps:
177             return pipe.named_steps[nm]
178
179     # 2) Por tipo
180     for nm, obj in pipe.named_steps.items():
181         if isinstance(obj, expected_type):
182             return obj
183
184     raise AssertionError(
185         f"No se ha encontrado en el Pipeline ningún paso con nombres {candidate_names}"
186         )
187         f"ni de tipo {expected_type.__name__}."
188
189 # Extraemos los bloques num y cat de forma robusta
190 num_trans, num_cols = _get_ct_block(preprocess, "num")
191 cat_trans, cat_cols = _get_ct_block(preprocess, "cat")
192
193 assert isinstance(num_trans, Pipeline), "✗ num_trans no es Pipeline."
194 assert isinstance(cat_trans, Pipeline), "✗ cat_trans no es Pipeline."
195
196 print("NUM cols:", num_cols)
197 print("CAT cols:", cat_cols)
198
199 # OJO: aquí usamos nombres reales ("imputer", "scaler", "ohe") pero con fallback por tipo
200 num_imp = _find_step(num_trans, ["imp", "imputer"], SimpleImputer)
201 num_sc = _find_step(num_trans, ["sc", "scaler"], StandardScaler)
202 cat_imp = _find_step(cat_trans, ["imp", "imputer"], SimpleImputer)
203 cat_oh = _find_step(cat_trans, ["oh", "ohe"], OneHotEncoder)
204
205 for obj, name in [(num_imp, "num_imputer"), (num_sc, "num_scaler"),
206                   (cat_imp, "cat_imputer"), (cat_oh, "cat_ohe")]:
207     assert obj is not None, f"✗ Falta paso {name} en preprocess."
208
209 print("✓ Estructura del preprocesador PRE compatible con export.")

```

```

210
211 """Celda 4 - export del preprocesador"""
212
213 # Celda 4 - Export portable del preprocesador PRE + centroides
214
215 def export_preprocess_pre_portable(preprocess, centroides, preprocess_meta, out_dir: Path):
216     # Usamos los helpers definidos en Celda 3
217     num_trans, num_cols = _get_ct_block(preprocess, "num")
218     cat_trans, cat_cols = _get_ct_block(preprocess, "cat")
219
220     # Localizamos los pasos clave de forma robusta
221     num_imp = _find_step(num_trans, ["imp", "imputer"], SimpleImputer)
222     num_sc = _find_step(num_trans, ["sc", "scaler"], StandardScaler)
223     cat_imp = _find_step(cat_trans, ["imp", "imputer"], SimpleImputer)
224     cat_oh = _find_step(cat_trans, ["oh", "ohe"], OneHotEncoder)
225
226     oh_categories = {}
227     for col, cats in zip(cat_cols, cat_oh.categories_):
228         oh_categories[col] = [str(c) for c in cats]
229
230     try:
231         feat_out = preprocess.get_feature_names_out()
232         feat_out = [str(f) for f in feat_out]
233     except Exception:
234         feat_out = None
235
236     raw_contract = (
237         preprocess_meta.get("expected_raw_features")
238         or preprocess_meta.get("expects_raw_features")
239         or {}
240     )
241
242     artifact = {
243         "artifact": "preprocess_pipeline.pkl",
244         "version": 1,
245         "created_at": datetime.datetime.utcnow().isoformat() + "Z",
246         "expected_raw_features": raw_contract,
247         "numeric": {
248             "columns": num_cols,
249             "imputer": {
250                 "strategy": num_imp.strategy,
251                 "statistics": _to_serializable(num_imp.statistics_)
252             },
253             "scaler": {
254                 "with_mean": bool(getattr(num_sc, "with_mean", True)),
255                 "with_std": bool(getattr(num_sc, "with_std", True)),
256                 "mean": _to_serializable(getattr(num_sc, "mean_", np.zeros(len(num_cols)))),
257                 "scale": _to_serializable(getattr(num_sc, "scale_", np.ones(len(num_cols))))
258             }
259         },
260         "categorical": {
261             "columns": cat_cols,
262             "imputer": {
263                 "strategy": cat_imp.strategy,
264                 "fill_values": _to_serializable(cat_imp.statistics_)
265             },
266             "onehot": {
267                 "handle_unknown": getattr(cat_oh, "handle_unknown", "ignore"),
268                 "categories": oh_categories
269             }
270         },
271         "feature_names_out": feat_out,
272         "centroids_ref": "centroides_pre.json",
273         "centroids": {
274             "classes": [str(c) for c in centroides.get("classes", [])],
275             "feature_names": centroides.get("feature_names"),
276             "centroids_pre": {str(k): _to_serializable(v) for k, v in centroides.get("centroids_pre", {}).items()},
277             "p95_intra": {str(k): float(v) for k, v in centroides.get("p95_intra", )}
278         }
279     }

```

```

        {}).items()
    }
}

out_path = out_dir / "preprocess_pre_portable.json"
save_json(artifact, out_path)
return out_path

pre_pre_path = export_preprocess_pre_portable(preprocess, centroides, preprocess_meta,
OUT_DIR)

"""Celda 5 - Export portable de RandomForest (RF base y RF_124)"""

# Celda 5 - Export portable de RandomForest (RF base y RF_124)

from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import _tree
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler, OneHotEncoder

def _extract_prep_from_pipe(pipe: Pipeline):
    """
    Extrae esquema de preprocesado (ColumnTransformer 'prep') de un Pipeline sklearn
    y el clasificador RandomForest ('clf'), en formato portable.
    """
    assert isinstance(pipe, Pipeline), "pipe debe ser sklearn.pipeline.Pipeline"
    prep = pipe.named_steps.get("prep")
    clf = pipe.named_steps.get("clf")
    assert isinstance(prep, ColumnTransformer), "X paso 'prep' no es
    ColumnTransformer."
    assert isinstance(clf, RandomForestClassifier), "X paso 'clf' no es
    RandomForestClassifier."

    # Localizamos bloques 'num' y 'cat' sin usar dict(...)
    num_trans = cat_trans = None
    num_cols = cat_cols = None
    for name, trans, cols in prep.transformers_:
        if name == "num":
            num_trans = trans
            num_cols = list(cols)
        elif name == "cat":
            cat_trans = trans
            cat_cols = list(cols)

    assert num_trans is not None and num_cols is not None, "X Bloque 'num' no
    encontrado en prep."
    assert cat_trans is not None and cat_cols is not None, "X Bloque 'cat' no
    encontrado en prep."
    assert isinstance(num_trans, Pipeline), "X num_trans no es Pipeline."
    assert isinstance(cat_trans, Pipeline), "X cat_trans no es Pipeline."

    # Usamos _find_step para no depender de nombres exactos ("imputer"/"scaler"/"ohe")
    num_imp = _find_step(num_trans, ["imp", "imputer"], SimpleImputer)
    num_sc = _find_step(num_trans, ["sc", "scaler"], StandardScaler)
    cat_imp = _find_step(cat_trans, ["imp", "imputer"], SimpleImputer)
    cat_oh = _find_step(cat_trans, ["oh", "ohe"], OneHotEncoder)

    oh_categories = {}
    for col, cats in zip(cat_cols, cat_oh.categories_):
        oh_categories[col] = [str(c) for c in cats]

    try:
        feat_out = prep.get_feature_names_out()
        feat_out = [str(f) for f in feat_out]
    except Exception:
        feat_out = None

    prep_schema = {
        "numeric": {
            "columns": num_cols,

```

```

344
345     "imputer": {
346         "strategy": num_imp.strategy,
347         "statistics": _to_serializable(num_imp.statistics_)
348     },
349     "scaler": {
350         "with_mean": bool(getattr(num_sc, "with_mean", True)),
351         "with_std": bool(getattr(num_sc, "with_std", True)),
352         "mean": _to_serializable(getattr(num_sc, "mean_", np.zeros(len(
353             num_cols)))), "scale": _to_serializable(getattr(num_sc, "scale_", np.ones(len(
354             num_cols))))
355     },
356     "categorical": {
357         "columns": cat_cols,
358         "imputer": {
359             "strategy": cat_imp.strategy,
360             "fill_values": _to_serializable(cat_imp.statistics_)
361         },
362         "onehot": {
363             "handle_unknown": getattr(cat_oh, "handle_unknown", "ignore"),
364             "categories": oh_categories
365         }
366     },
367     "feature_names_out": feat_out,
368 }
369
370     return prep_schema, clf
371
372 def _tree_to_dict(tree: _tree.Tree):
373     """Convierte un sklearn.tree._tree.Tree en un dict portable."""
374     n_nodes = tree.node_count
375     return {
376         "n_nodes": int(n_nodes),
377         "children_left": _to_serializable(tree.children_left),
378         "children_right": _to_serializable(tree.children_right),
379         "feature": _to_serializable(tree.feature),
380         "threshold": _to_serializable(tree.threshold),
381         "value": _to_serializable(tree.value.squeeze(axis=1)), # [n_nodes,
382         n_classes]
383         "impurity": _to_serializable(tree.impurity),
384         "n_node_samples": _to_serializable(tree.n_node_samples)
385     }
386
387 def export_rf_portable(pipe: Pipeline, meta: dict, out_name: str, out_dir: Path):
388     """
389     Exporta un Pipeline (prep+clf) de RF a JSON portable:
390     - feature_space.expected_raw_features
391     - feature_space.alias_mapping_used_in_train
392     - feature_space.dtypes_expect
393     - esquema de preprocessado (prep_schema)
394     - bosque {trees: [...]} árbol a árbol
395     """
396     prep_schema, clf = _extract_prep_from_pipe(pipe)
397
398     trees = []
399     for est in clf.estimators_:
400         trees.append(_tree_to_dict(est.tree_))
401
402     classes = [str(c) for c in clf.classes_]
403
404     feature_space = {
405         "expected_raw_features": meta.get("expected_raw_features", {}),
406         "alias_mapping_used_in_train": meta.get("alias_mapping_used_in_train", {}),
407         "dtypes_expect": meta.get("dtypes_expect", {}),
408         "preprocess": prep_schema
409     }
410
411     artifact = {
412         "artifact": meta.get("artifact", out_name.replace(".json", ".pkl")),
413         "model_type": "RandomForestClassifier",
414         "version": 1,
415     }

```

```

413         "created_at": datetime.datetime.utcnow().isoformat() + "Z",
414         "classes": classes,
415         "n_classes": int(len(classes)),
416         "feature_space": feature_space,
417         "forest": {
418             "n_estimators": int(clf.n_estimators),
419             "max_depth": int(clf.max_depth) if clf.max_depth is not None else None,
420             "n_features_in": int(getattr(clf, "n_features_in_", 0)),
421             "trees": trees
422         }
423     }
424
425     out_path = out_dir / out_name
426     save_json(artifact, out_path)
427     return out_path
428
429 # Export de RF base y RF_124
430 rf_base_portable_path = export_rf_portable(rf_base, rf_base_meta,
431 "rf_base_portable.json", OUT_DIR)
431 rf_124_portable_path = export_rf_portable(rf_124, rf_124_meta, "rf_124_portable.json"
432 , OUT_DIR)
432 print("✓ Export RF base:", rf_base_portable_path)
433 print("✓ Export RF 124 :", rf_124_portable_path)
434
435 """Celda 6 - Export portable de XGBoost (baja_comp)"""
436
437 # Celda 6 - Export portable de XGBoost (baja_comp)
438
439 if not HAS_XGB:
440     print("⚠️ XGBoost no disponible; sólo se exportarán metadatos sin booster JSON.")
441 else:
442     # Modelo entrenado y meta original
443     xgb_model = xgb_pack["modelo"]
444     xgb_vars = xgb_meta.get("variables_predictoras", [])
445     enc_cols = xgb_meta.get("encoders_columns", [])
446     labels = xgb_meta.get("labels", [])
447     policies = xgb_meta.get("policies", {})
448     threshold_bajo = policies.get("threshold_bajo", None)
449
450     # Encoders usados en train (LabelEncoder/OrdinalEncoder normalmente)
451     label_encoders = xgb_pack.get("encoders", {})
452
453     encoders_schema = {}
454     for col, le in label_encoders.items():
455         classes = getattr(le, "classes_", None)
456         encoders_schema[col] = {
457             "classes": [str(c) for c in classes] if classes is not None else []
458         }
459
460     # Meta portable = copia enriquecida del meta original
461     xgb_meta_portable = dict(xgb_meta)
462     xgb_meta_portable["created_at"] = datetime.datetime.utcnow().isoformat() + "Z"
463     xgb_meta_portable["artifact_local"] = "model_Noferta_XGB_v8.pkl"
464     xgb_meta_portable["encoders_schema"] = encoders_schema
465
466     meta_path = OUT_DIR / "xgb_baja_portable.json"
467     save_json(xgb_meta_portable, meta_path)
468
469     # Booster JSON (modelo en formato nativo XGBoost)
470     booster = xgb_model.get_booster()
471     booster_path = OUT_DIR / "xgb_baja_model.json"
472     booster.save_model(str(booster_path))
473     print(f"▣ Guardado booster XGB en {booster_path}")
474     print(f"▣ Guardado meta portable XGB en {meta_path}")
475
476 """Celda 7 - Export portable de kNN (espacio PRE_NUM + matrices)"""
477
478 # Celda 7 - Export portable de kNN (espacio PRE_NUM + matrices)
479
480 from sklearn.neighbors import NearestNeighbors
481 from sklearn.impute import SimpleImputer
482 from sklearn.preprocessing import StandardScaler

```

```

483
484     assert isinstance(nn_train, NearestNeighbors), "X nn_train no es NearestNeighbors."
485     assert isinstance(nn_all, NearestNeighbors), "X nn_all no es NearestNeighbors."
486
487     Z_train = getattr(nn_train, "_fit_X", None)
488     Z_all = getattr(nn_all, "_fit_X", None)
489     assert Z_train is not None, "X nn_train no tiene _fit_X."
490     assert Z_all is not None, "X nn_all no tiene _fit_X."
491
492     np.savez_compressed(OUT_DIR / "knn_train.npz", Z=Z_train)
493     np.savez_compressed(OUT_DIR / "knn_all.npz", Z=Z_all)
494     print("Guardadas matrices knn_train.npz y knn_all.npz")
495
496     # Extraer esquema de num_pipeline de forma robusta
497     num_imp = _find_step(knn_num_pipe, ["imp", "imputer"], SimpleImputer)
498     num_sc = _find_step(knn_num_pipe, ["sc", "scaler"], StandardScaler)
499
500     knn_space_portable = {
501         "artifact_train": "knn_pre_train.joblib",
502         "artifact_all": "knn_pre_all.joblib",
503         "version": 1,
504         "created_at": datetime.datetime.utcnow().isoformat() + "Z",
505         "pre_num_features": knn_params.get("pre_num_features", []),
506         "metric": knn_params.get("metric", getattr(nn_all, "metric", "euclidean")),
507         "k_neighbors": int(knn_params.get("k_neighbors", getattr(nn_all, "n_neighbors", 25))),
508         "p95_nn1_train": float(knn_params.get("p95_nn1_train", 0.0)),
509         "num_pipeline": {
510             "imputer": {
511                 "strategy": num_imp.strategy,
512                 "statistics": _to_serializable(num_imp.statistics_)
513             },
514             "scaler": {
515                 "with_mean": bool(getattr(num_sc, "with_mean", True)),
516                 "with_std": bool(getattr(num_sc, "with_std", True)),
517                 "mean": _to_serializable(getattr(num_sc, "mean_", np.zeros(Z_train.shape[1]))),
518                 "scale": _to_serializable(getattr(num_sc, "scale_", np.ones(Z_train.shape[1])))
519             }
520         },
521         "catalog": {
522             "id_col": knn_id_col,
523             "index_key": knn_index_key,
524             "source_df": knn_catalog_meta.get("source_df") if knn_catalog_meta else None,
525             "rows": knn_catalog_meta.get("rows") if knn_catalog_meta else None,
526             "cols": knn_catalog_meta.get("cols") if knn_catalog_meta else None,
527         }
528     }
529
530     save_json(knn_space_portable, OUT_DIR / "knn_space_portable.json")
531     print("✓ Export KNN portable completo.")
532
533     """Celda 8 - Export portable de ensemble"""
534
535     # Celda 8 - Export portable de ensemble
536
537     ensemble_portable = dict(ensemble_meta)  # copia superficial del meta original
538
539     ensemble_portable["version"] = 1
540     ensemble_portable["created_at"] = datetime.datetime.utcnow().isoformat() + "Z"
541
542     # Añadimos referencias explícitas a los artefactos portables
543     depends_on_portable = {
544         "rf_base": "rf_base_portable.json",
545         "rf_124_basebc": "rf_124_portable.json",
546         "xgb_baja": "xgb_baja_portable.json" if HAS_XGB else None,
547         "preprocess_pre": "preprocess_pre_portable.json",
548         "knn_space": "knn_space_portable.json",
549         "centroids": "centroides_pre.json",
550     }
551

```

```

552 ensemble_portable["depends_on_portable"] = depends_on_portable
553
554 save_json(ensemble_portable, OUT_DIR / "ensemble_portable.json")
555 print("✓ Export ensemble_portable.json")
556
557 """Celda 9 - Test de equivalencia completo para RF base (portable vs original)"""
558
559 # Celda 9 - Test de equivalencia completo para RF base (portable vs original, FIX
560 # float32)
561
562 from collections import OrderedDict
563 import numpy as np
564 import pandas as pd
565 import datetime
566
567 # -----
568 # 1. Cargamos datos de test
569 # -----
570 df_full = pd.read_excel(TEST_XLSX)
571 assert "Dataset" in df_full.columns, "✗ Falta columna Dataset en df_test_train_v8.xlsx."
572
573 df_test = df_full[df_full["Dataset"].astype(str).str.upper() == "TEST"].copy()
574 df_test.reset_index(drop=True, inplace=True)
575
576 print(f"Filas en TEST: {len(df_test)}")
577
578 # -----
579 # 2. Cargamos modelos portables RF
580 # -----
581 rf_base_port = json.loads((OUT_DIR / "rf_base_portable.json").read_text(encoding="utf-8"))
582 rf_124_port = json.loads((OUT_DIR / "rf_124_portable.json").read_text(encoding="utf-8"))
583
584 rf_base_clf = rf_base.named_steps["clf"]
585 rf_124_clf = rf_124.named_steps["clf"]
586
587 # -----
588 # 3. Helpers de inferencia portable (preprocesado + árboles)
589 # -----
590
591 def _portable_preprocess_df(df_raw: pd.DataFrame, prep_schema: dict) -> np.ndarray:
592     """
593         Aplica el preprocesado descrito en prep_schema (tal y como se exportó del
594         ColumnTransformer):
595             - numeric: imputación + escalado
596             - categorical: imputación + one-hot
597
598         En categóricas detecta automáticamente si las categorías son numéricas
599         (caso 'Mes_lici' = 1..12) para evitar errores tipo "1" vs "1.0".
600     """
601
602     # --- Numérico ---
603     num_conf = prep_schema["numeric"]
604     num_cols = list(num_conf["columns"])
605     if num_cols:
606         X_num = df_raw[num_cols].apply(pd.to_numeric, errors="coerce")
607         X_num = X_num.to_numpy(dtype=float)
608
609         stats = np.array(num_conf["imputer"]["statistics"], dtype=float)
610         assert stats.shape[0] == X_num.shape[1], "Dimensión de estadísticas numéricas
611         no cuadra."
612
613         for j in range(X_num.shape[1]):
614             mask = np.isnan(X_num[:, j])
615             if mask.any():
616                 X_num[mask, j] = stats[j]
617
618             sc_conf = num_conf["scaler"]
619             with_mean = bool(sc_conf.get("with_mean", True))
620             with_std = bool(sc_conf.get("with_std", True))
621             mean = np.array(sc_conf.get("mean", np.zeros_like(stats)), dtype=float)

```

```

618     scale = np.array(sc_conf.get("scale", np.ones_like(stats)), dtype=float)
619
620     Xn = X_num.copy()
621     if with_mean:
622         Xn = Xn - mean
623     if with_std:
624         scale_safe = np.where(scale == 0, 1.0, scale)
625         Xn = Xn / scale_safe
626 else:
627     Xn = np.zeros((len(df_raw), 0), dtype=float)
628
629 # --- Categórico ---
630 cat_conf = prep_schema["categorical"]
631 cat_cols = list(cat_conf["columns"])
632 if cat_cols:
633     fill_values = list(cat_conf["imputer"]["fill_values"])
634     assert len(fill_values) == len(cat_cols), "Dimensión de fill_values no cuadra con columnas categóricas."
635
636     cat_onehot_conf = cat_conf["onehot"]
637     categories_map = cat_onehot_conf["categories"] # dict {col: [cats...]}
638
639     Xc_blocks = []
640     for j, col in enumerate(cat_cols):
641         s = df_raw[col].copy()
642
643         # Imputación "most_frequent" según estadísticas guardadas
644         fill_val = fill_values[j]
645         s = s.where(~s.isna(), other=fill_val)
646
647         cats = categories_map.get(col, [])
648         if not cats:
649             Xc_j = np.zeros((len(s), 0), dtype=float)
650             Xc_blocks.append(Xc_j)
651             continue
652
653         # Detectamos si todas las categorías son numéricas (caso Mes_lici, etc.)
654         numeric_cats = True
655         cats_num = []
656         for c in cats:
657             try:
658                 cats_num.append(float(c))
659             except (TypeError, ValueError):
660                 numeric_cats = False
661                 break
662
663         n = len(s)
664         if numeric_cats:
665             # Col categórica codificada como número: trabajamos en float y
666             # comparamos numéricamente
667             s_num = pd.to_numeric(s, errors="coerce")
668             s_arr = s_num.to_numpy(dtype=float)
669             cats_arr = np.array(cats_num, dtype=float)
670
671             Xc_j = np.zeros((n, len(cats_arr)), dtype=float)
672             for k, cat_val in enumerate(cats_arr):
673                 mask = (s_arr == cat_val)
674                 Xc_j[mask, k] = 1.0
675
676         else:
677             # Col categórica puramente textual: trabajamos en str
678             s_str = s.astype(str)
679             cats_str = [str(c) for c in cats]
680
681             Xc_j = np.zeros((n, len(cats_str)), dtype=float)
682             val_to_index = {v: idx for idx, v in enumerate(cats_str)}
683             for i, v in enumerate(s_str):
684                 idx = val_to_index.get(v, None)
685                 if idx is not None:
686                     Xc_j[i, idx] = 1.0
687
688             Xc_blocks.append(Xc_j)

```

```

688
689     Xc = np.concatenate(Xc_blocks, axis=1) if Xc_blocks else np.zeros((len(df_raw)
690         ), 0), dtype=float)
690
691     else:
692         Xc = np.zeros((len(df_raw), 0), dtype=float)
693
694 # Orden final: primero num, luego cat (igual que ColumnTransformer con bloques
695 # ["num", "cat"])
696 X_pre = np.concatenate([Xn, Xc], axis=1)
697 return X_pre
698
699
700 def _tree_predict_proba_single(X: np.ndarray, tree_dict: dict, n_classes: int) -> np.
701 ndarray:
702     """
703     Recorre un árbol individual (exportado) para todos los samples de X y devuelve
704     una matriz [n_samples, n_classes] con las probabilidades de ese árbol.
705
706     IMPORTANTE: emula el comportamiento de sklearn casteando a float32
707     antes del descenso, porque los árboles internamente comparan en float32.
708     """
709     children_left = np.array(tree_dict["children_left"], dtype=int)
710     children_right = np.array(tree_dict["children_right"], dtype=int)
711     feature = np.array(tree_dict["feature"], dtype=int)
712     threshold = np.array(tree_dict["threshold"], dtype=float)
713     value = np.array(tree_dict["value"], dtype=float) # [n_nodes, n_classes]
714
715     # Aquí está la clave para reproducir exactamente sklearn:
716     X32 = X.astype(np.float32, copy=False)
717     thr32 = threshold.astype(np.float32, copy=False)
718
719     n_samples = X32.shape[0]
720     proba = np.zeros((n_samples, n_classes), dtype=float)
721
722     for i in range(n_samples):
723         node = 0
724         # Descenso hasta hoja
725         while children_left[node] != -1: # _tree.TREE_LEAF == -1
726             f = feature[node]
727             t = thr32[node]
728             x_val = X32[i, f]
729             # Misma regla que sklearn, pero en float32
730             if x_val <= t:
731                 node = children_left[node]
732             else:
733                 node = children_right[node]
734
735             counts = value[node]
736             tot = counts.sum()
737             if tot <= 0:
738                 proba[i, :] = 1.0 / n_classes
739             else:
740                 proba[i, :] = counts / tot
741
742     return proba
743
744 def rf_portable_predict_proba(df_raw: pd.DataFrame, rf_port: dict) -> (np.ndarray, np.
745 ndarray):
746     """
747     Calcula predict_proba con el modelo RF portable:
748     - Reconstruye el preprocesado desde rf_port["feature_space"]["preprocess"].
749     - Recorre todos los árboles del bosque.
750     Devuelve (probas, classes) donde:
751     - probas: [n_samples, n_classes]
752     - classes: array de etiquetas en el orden de columnas de probas.
753     """
754     feature_space = rf_port["feature_space"]
755     prep_schema = feature_space["preprocess"]
756     expected = feature_space.get("expected_raw_features", {})
757     num_cols = expected.get("num", [])
758     cat_cols = expected.get("cat", [])
759     all_cols = list(num_cols) + list(cat_cols)

```

```

756     missing = [c for c in all_cols if c not in df_raw.columns]
757     assert not missing, f"✗ Faltan columnas esperadas en df_raw para RF portable: {missing}"
758
759     df_in = df_raw[all_cols].copy()
760
761     # Preprocesado portable
762     X_pre = _portable_preprocess_df(df_in, prep_schema)
763
764     forest = rf_port["forest"]
765     trees = forest["trees"]
766     n_estimators = int(forest["n_estimators"])
767     assert n_estimators == len(trees), "✗ n_estimators no cuadra con nº de árboles exportados."
768
769     classes = np.array(rf_port["classes"])
770     n_classes = int(rf_port["n_classes"])
771     assert n_classes == len(classes), "✗ n_classes no cuadra con longitud de classes."
772
773     proba_sum = np.zeros((X_pre.shape[0], n_classes), dtype=float)
774     for t_dict in trees:
775         proba_tree = _tree_predict_proba_single(X_pre, t_dict, n_classes)
776         proba_sum += proba_tree
777
778     proba = proba_sum / n_estimators
779     return proba, classes
780
781 def rf_portable_predict(df_raw: pd.DataFrame, rf_port: dict) -> np.ndarray:
782     proba, classes = rf_portable_predict_proba(df_raw, rf_port)
783     idx_max = np.argmax(proba, axis=1)
784     return classes[idx_max]
785
786 # -----
787 # 4. Equivalencia RF base (predict + predict_proba)
788 # -----
789
790 # Columnas esperadas por RF base según meta portable
791 fs_base = rf_base_port["feature_space"]
792 exp_base = fs_base.get("expected_raw_features", {})
793 num_base = exp_base.get("num", [])
794 cat_base = exp_base.get("cat", [])
795 cols_base = list(num_base) + list(cat_base)
796
797 missing_base = [c for c in cols_base if c not in df_test.columns]
798 if missing_base:
799     raise ValueError(f"✗ Faltan columnas requeridas por RF base en df_test: {missing_base}")
800
801 df_rf_base = df_test[cols_base].copy()
802
803 # Predicciones originales (Pipeline sklearn)
804 proba_orig_base = rf_base.predict_proba(df_rf_base)
805 pred_orig_base = rf_base.predict(df_rf_base)
806 classes_orig_base = rf_base_clf.classes_
807
808 # Predicciones portables (a partir del JSON)
809 proba_port_base, classes_port_base = rf_portable_predict_proba(df_rf_base,
810                     rf_base_port)
811 pred_port_base = rf_portable_predict(df_rf_base, rf_base_port)
812
813 # Alineamos posibles tipos en las labels
814 classes_orig_base_str = np.array([str(c) for c in classes_orig_base])
815 classes_port_base_str = np.array([str(c) for c in classes_port_base])
816
817 assert np.array_equal(classes_orig_base_str, classes_port_base_str), \
818     "✗ El orden/clases de RF base original vs portable no coincide."
819
820 # Comprobaciones de equivalencia
821 pred_equal_base = np.array_equal(
822     np.array([str(p) for p in pred_orig_base]),
823     np.array([str(p) for p in pred_port_base]))

```

```

823 )
824
825 max_abs_diff_base = float(np.max(np.abs(proba_orig_base - proba_port_base)))
826
827 # Tolerancia muy estricta (debería ser prácticamente 0 tras el FIX)
828 PROBA_TOL = 1e-12
829 proba_close_base = bool(max_abs_diff_base <= PROBA_TOL)
830
831 print(f"RF base - pred_equal: {pred_equal_base}, proba_max_abs_diff: {max_abs_diff_base:.3e}")
832
833 rf_base_strict_equiv_ok = bool(pred_equal_base and proba_close_base)
834
835 # -----
836 # 5. Metadatos de coherencia y placeholders para el resto
837 # -----
838 report = OrderedDict()
839 report["created_at"] = datetime.datetime.utcnow().isoformat() + "Z"
840 report["n_test_rows"] = int(len(df_test))
841
842 # Metadatos RF base / RF 124
843 report["rf_base_meta_ok"] = (
844     len(rf_base_port["forest"]["trees"]) == int(getattr(rf_base_clf, "n_estimators", -1))
845     and sorted(rf_base_port["classes"]) == sorted([str(c) for c in rf_base_clf.classes_])
846 )
847
848 report["rf_124_meta_ok"] = (
849     len(rf_124_port["forest"]["trees"]) == int(getattr(rf_124_clf, "n_estimators", -1))
850     and sorted(rf_124_port["classes"]) == sorted([str(c) for c in rf_124_clf.classes_])
851 )
852
853 # XGB: comprobamos sólo que el booster portable se puede cargar
854 if HAS_XGB:
855     try:
856         booster_loaded = xgb.Booster()
857         booster_loaded.load_model(str(OUT_DIR/"xgb_baja_model.json"))
858         report["xgb_booster_load_ok"] = True
859     except Exception as e:
860         report["xgb_booster_load_ok"] = False
861         report["xgb_booster_load_error"] = str(e)
862     else:
863         report["xgb_booster_load_ok"] = None
864
865 # Resultados detallados RF base
866 report["rf_base_strict_equiv_ok"] = rf_base_strict_equiv_ok
867 report["rf_base_pred_equal"] = bool(pred_equal_base)
868 report["rf_base_proba_max_abs_diff"] = max_abs_diff_base
869 report["rf_base_proba_tol"] = PROBA_TOL
870
871 # RF 124 y pipeline completo -> pendientes de implementar (de momento)
872 report["rf_124_strict_equiv_ok"] = None
873 report["full_pipeline_equiv_ok"] = None
874
875 # Criterio de OK global (de momento sólo RF base + coherencia básica)
876 report["ok"] = bool(
877     report["rf_base_meta_ok"]
878     and rf_base_strict_equiv_ok
879     and (report["xgb_booster_load_ok"] in [True, None])
880 )
881
882 equiv_path = OUT_DIR / "equivalence_report.json"
883 save_json(report, equiv_path)
884 print("▣ Equivalence report actualizado en:", equiv_path)
885 print("✓ OK global (parcial, centrado en RF base):", report["ok"])

```