

```

1  from __future__ import annotations
2
3  ### SCRIPT 19 - PYTHON
4  # PIPELINE EJECUCION MODELOS EN GPT
5  # PIPELINE PREDICCION MODELOS DESDE ARCHIVOS JSON y NPZ
6  #
7  # =====
8
9 """
10 Módulo de inferencia portable para el pipeline de riesgo en licitaciones.
11
12 - Lee artefactos JSON/NPZ desde un ZIP (modelos_portable.zip).
13 - Implementa versiones "puras" en NumPy/Pandas de:
14   - Preprocesadores num+cat con imputación + escalado + one-hot.
15   - RandomForestClassifier (RF base y RF_124_BASE+BC).
16   - XGBoost multi-clase (BAJO/MEDIO/ALTO) para baja_comp.
17   - kNN en el espacio PRE_NUM.
18   - Ensemble RF base vs RF_124 con reglas simplificadas.
19
20 Notas importantes:
21 - No depende de scikit-learn ni xgboost.
22 - Asume una tolerancia práctica de ≈0,02 en predict_proba de RF
23   respecto a sklearn; predict (las clases) debe coincidir si se respeta
24   el contrato de entrada.
25 """
26
27 import json
28 import math
29 import zipfile
30 from dataclasses import dataclass
31 from typing import Any, Dict, List, Optional, Tuple
32
33 import numpy as np
34 import pandas as pd
35
36
37 # -----
38 # Utilidades numéricas básicas
39 # -----
40
41 def _softmax(x: np.ndarray) -> np.ndarray:
42     """
43         Softmax fila a fila sobre un array 2D.
44
45         Si se pasa un vector 1D, se convierte internamente en 2D (shape (1, n)).
46     """
47     x = np.asarray(x, dtype=float)
48     if x.ndim == 1:
49         x = x.reshape(1, -1)
50     x_max = np.max(x, axis=1, keepdims=True)
51     e = np.exp(x - x_max)
52     s = e.sum(axis=1, keepdims=True)
53     # Evitar divisiones por cero
54     s[s == 0.0] = 1.0
55     return e / s
56
57
58 def _ensure_2d_array(X: Any) -> np.ndarray:
59     """
60         Asegura que la entrada se convierte en un array 2D de tipo float.
61
62         - Si es un DataFrame, usa .to_numpy().
63         - Si es un vector 1D, lo convierte en shape (1, n_features).
64     """
65     if isinstance(X, pd.DataFrame):
66         A = X.to_numpy(dtype=float)
67     else:
68         A = np.asarray(X, dtype=float)
69     if A.ndim == 1:
70         A = A.reshape(1, -1)
71     return A
72

```

```

73
74 def _is_int_string(s: str) -> bool:
75     """Devuelve True si la cadena representa un entero (signo opcional)."""
76     if s is None:
77         return False
78     s = str(s).strip()
79     if not s:
80         return False
81     if s[0] in "+-":
82         s = s[1:]
83     return s.isdigit()
84
85 """Celda 2- ZipStore"""
86
87 # -----
88 # ZipStore: lectura cómoda de JSON/NPZ desde un ZIP
89 # -----
90
91 @dataclass
92 class ZipStore:
93     """
94     Pequeño contenedor para acceder a artefactos dentro de un ZIP.
95
96     Uso típico:
97         store = ZipStore("/ruta/a/modelos_portable.zip")
98         meta_rf = store.read_json("rf_base_portable.json")
99         Z_train = store.read_npz("knn_train.npz")["Z"]
100
101    Internamente construye un índice de:
102        nombre_corto -> ruta_completa_en_el_zip
103
104    por ejemplo:
105        "rf_base_portable.json" -> "modelos_portable/rf_base_portable.json"
106    """
107
108    zip_path: str
109
110    def __post_init__(self) -> None:
111        # Abrimos el ZIP una única vez y construimos un índice por nombre corto
112        self._zf = zipfile.ZipFile(self.zip_path, mode="r")
113        self._index: Dict[str, str] = {}
114
115        for name in self._zf.namelist():
116            # Ignoramos directorios
117            if name.endswith("/"):
118                continue
119            short = name.split("/")[-1]
120            if short in self._index:
121                # Preferimos fallar pronto si hay colisiones de nombres cortos
122                raise ValueError(
123                    f"Nombre corto duplicado en el ZIP: {short!r} "
124                    f"(rutas: {self._index[short]!r} y {name!r})"
125                )
126            self._index[short] = name
127
128    # ----- Métodos públicos de utilidad -----
129
130    def list(self) -> List[str]:
131        """Devuelve la lista de nombres cortos disponibles en el ZIP."""
132        return sorted(self._index.keys())
133
134    def has(self, short_name: str) -> bool:
135        """Indica si existe un fichero con ese nombre corto en el ZIP."""
136        return short_name in self._index
137
138    def full_name(self, short_name: str) -> str:
139        """Resuelve el nombre corto a la ruta completa dentro del ZIP."""
140        try:
141            return self._index[short_name]
142        except KeyError:
143            raise FileNotFoundError(
144                f"No se encuentra {short_name!r} en el ZIP {self.zip_path!r}."

```

```

145
146
147     def read_json(self, short_name: str) -> Dict[str, Any]:
148         """
149             Lee y decodifica un JSON dentro del ZIP.
150
151             Parámetros
152             -----
153             short_name : str
154                 Nombre corto del fichero, por ejemplo 'rf_base_portable.json'.
155
156             Devuelve
157             -----
158             dict
159                 Contenido del JSON parseado.
160             """
161             full = self._full_name(short_name)
162             raw = self._zf.read(full)
163             return json.loads(raw.decode("utf-8"))
164
165     def read_npz(self, short_name: str) -> np.lib.npyio.NpzFile:
166         """
167             Lee un fichero .npz dentro del ZIP y devuelve el objeto NpzFile.
168
169             Parámetros
170             -----
171             short_name : str
172                 Nombre corto del fichero, por ejemplo 'knn_train.npz'.
173
174             Devuelve
175             -----
176             np.lib.npyio.NpzFile
177                 Objeto desde el que se pueden leer arrays, p.ej. ['Z'].
178             """
179             from io import BytesIO
180
181             full = self._full_name(short_name)
182             raw = self._zf.read(full)
183             bio = BytesIO(raw)
184             return np.load(bio)
185
186     def close(self) -> None:
187         """Cierra el ZIP subyacente (opcional; no es estrictamente necesario)."""
188         try:
189             self._zf.close()
190         except Exception:
191             pass
192
193         # Soporte para usarlo como contexto: with ZipStore(...) as store:
194     def __enter__(self) -> "ZipStore":
195         return self
196
197     def __exit__(self, exc_type, exc_val, exc_tb) -> None:
198         self.close()
199
200     """Celda 3 - Preprocesador portable (PortablePreprocessor)"""
201
202     # -----
203     # PortablePreprocessor
204     # -----
205
206     class PortablePreprocessor:
207         """
208             Preprocesador portable que reproduce:
209             - SimpleImputer (num) + StandardScaler
210             - SimpleImputer (cat) + OneHotEncoder(handle_unknown="ignore")
211
212             El esquema debe tener la forma usada en:
213             - preprocess_pre_portable.json
214             - feature_space["preprocess"] de los RF portables.
215         """

```

```

217     def __init__(self, schema: Dict[str, Any]) -> None:
218         self.schema = schema
219
220         num = schema.get("numeric", {}) or {}
221         cat = schema.get("categorical", {}) or {}
222
223         # ----- Configuración numérica -----
224         self.num_cols: List[str] = list(num.get("columns", []))
225         num_imp = num.get("imputer", {}) or {}
226         num_scl = num.get("scaler", {}) or {}
227
228         self.num_imp_strategy: str = num_imp.get("strategy", "median")
229         self.num_imp_statistics: np.ndarray = np.asarray(
230             num_imp.get("statistics", []), dtype=float
231         )
232
233         self.num_with_mean: bool = bool(num_scl.get("with_mean", True))
234         self.num_with_std: bool = bool(num_scl.get("with_std", True))
235         self.num_mean: np.ndarray = np.asarray(
236             num_scl.get("mean", np.zeros(len(self.num_cols))), dtype=float
237         )
238         self.num_scale: np.ndarray = np.asarray(
239             num_scl.get("scale", np.ones(len(self.num_cols))), dtype=float
240         )
241         # Evitar divisiones por cero más adelante
242         self.num_scale[self.num_scale == 0.0] = 1.0
243
244         # ----- Configuración categórica -----
245         self.cat_cols: List[str] = list(cat.get("columns", []))
246         cat_imp = cat.get("imputer", {}) or {}
247         cat_oh = cat.get("onehot", {}) or {}
248
249         self.cat_imp_strategy: str = cat_imp.get("strategy", "most_frequent")
250         self.cat_fill_values: List[Any] = list(cat_imp.get("fill_values", []))
251         self.cat_handle_unknown: str = cat_oh.get("handle_unknown", "ignore")
252
253         # categories: dict col -> list[str]
254         self.cat_categories: Dict[str, List[str]] = cat_oh.get("categories", {}) or {}
255
256         # Nombres de salida opcionales (no imprescindibles para la inferencia)
257         self.feature_names_out: Optional[List[str]] = schema.get("feature_names_out")
258
259         # ----- numéricos -----
260
261     def _transform_numeric(self, df: pd.DataFrame) -> np.ndarray:
262         if not self.num_cols:
263             return np.zeros((len(df), 0), dtype=float)
264
265         X_num = pd.DataFrame(index=df.index)
266         for c in self.num_cols:
267             if c in df.columns:
268                 X_num[c] = pd.to_numeric(df[c], errors="coerce")
269             else:
270                 # Columna ausente: todo NaN
271                 X_num[c] = np.nan
272
273         A = X_num.to_numpy(dtype=float) # shape (n_samples, n_features)
274
275         # Imputación: implementamos 'median'/'most_frequent' vía statistics
276         if self.num_imp_strategy in ("median", "most_frequent"):
277             stats = self.num_imp_statistics
278             if stats.shape[0] != A.shape[1]:
279                 raise ValueError(
280                     "Dimensión de estadísticas numéricas inconsistente: "
281                     f"{stats.shape[0]} != {A.shape[1]}"
282                 )
283             # Imputamos NaN columna a columna
284             for j in range(A.shape[1]):
285                 col = A[:, j]
286                 mask = np.isnan(col)
287                 if mask.any():
288                     col[mask] = stats[j]

```

```

289     # A ya modificado in-place
290
291     # Escalado estándar
292     if self.num_with_mean:
293         A = A - self.num_mean
294     if self.num_with_std:
295         A = A / self.num_scale
296
297     return A
298
299     # ----- categóricas -----
300
301 def _normalize_categorical_values(self, col: str, s: pd.Series) -> pd.Series:
302     """
303     Normaliza los valores de una columna categórica:
304
305     - Si todas las categorías de esquema son enteros representados como str
306       ("1", "2", ...), convierte los valores numéricos a int y luego a str
307       (así 1, 1.0 y "1" se mapean a "1").
308     - En otro caso, convierte a str tal cual.
309     """
310
311     cats = self.cat_categories.get(col, [])
312     if cats and all(_is_int_string(c) for c in cats):
313         # Columna con categorías "numéricas"
314         def to_int_str(x: Any) -> Any:
315             if pd.isna(x):
316                 return np.nan # se imputará antes
317             sx = str(x).strip()
318             try:
319                 # Permite valores tipo "1.0"
320                 val = int(float(sx))
321                 return str(val)
322             except Exception:
323                 return sx # fallback: lo que sea como str
324
325         return s.map(to_int_str)
326     else:
327         # Columna categórica convencional
328         return s.astype(str)
329
330 def _transform_categorical(self, df: pd.DataFrame) -> np.ndarray:
331     if not self.cat_cols:
332         return np.zeros((len(df), 0), dtype=float)
333
334     n = len(df)
335     blocks: List[np.ndarray] = []
336
337     for idx, col in enumerate(self.cat_cols):
338         if col in df.columns:
339             s = df[col]
340         else:
341             # Columna ausente: todo NaN
342             s = pd.Series([np.nan] * n, index=df.index)
343
344         # Imputación
345         if self.cat_imp_strategy == "most_frequent":
346             # fill_values está en el mismo orden que cat_cols
347             if idx < len(self.cat_fill_values):
348                 fill_val = self.cat_fill_values[idx]
349             else:
350                 # Fallback defensivo
351                 fill_val = (
352                     s.mode(dropna=True).iloc[0] if s.notna().any() else ""
353                 )
354             s = s.fillna(fill_val)
355         else:
356             # Otras estrategias no se contemplan en tus artefactos
357             s = s.fillna("")
358
359         # Normalización de valores (numéricos como enteros en str)
360         s_norm = self._normalize_categorical_values(col, s)

```

```

361     # One-hot
362     cats = self.cat_categories.get(col, [])
363     cats = [str(c) for c in cats]
364
365     # Matriz (n_samples, len(cats))
366     M = np.zeros((n, len(cats)), dtype=float)
367     if len(cats) == 0:
368         blocks.append(M)
369         continue
370
371     # Mapeo valor->índice de categoría
372     cat_to_idx = {str(c): j for j, c in enumerate(cats)}
373
374     # Para cada fila, activamos la posición correspondiente
375     for i, val in enumerate(s_norm):
376         if pd.isna(val):
377             # Ya imputado, pero por seguridad
378             continue
379         key = str(val)
380         j = cat_to_idx.get(key, None)
381         if j is None:
382             # handle_unknown:
383             if self.cat_handle_unknown == "ignore":
384                 continue
385             else:
386                 # En tus modelos, esto no debería ocurrir; lanzamos error
387                 # explícito
388                 raise ValueError(
389                     f"Valor desconocido '{key}' en columna '{col}' "
390                     f"con handle_unknown='{self.cat_handle_unknown}'"
391                 )
392         M[i, j] = 1.0
393
394     blocks.append(M)
395
396     if blocks:
397         return np.concatenate(blocks, axis=1)
398     else:
399         return np.zeros((len(df), 0), dtype=float)
400
401     # ----- API pública -----
402
403     def transform(self, df: pd.DataFrame) -> np.ndarray:
404         """
405             Aplica el preprocessado numérico + categórico y devuelve
406             una matriz NumPy 2D (n_samples, n_features_transformados).
407         """
408         if not isinstance(df, pd.DataFrame):
409             raise TypeError("PortablePreprocessor.transform espera un
410                             pandas.DataFrame.")
411
412         A_num = self._transform_numeric(df)
413         A_cat = self._transform_categorical(df)
414
415         if A_num.shape[0] != A_cat.shape[0]:
416             raise ValueError("Dimensiones numéricas y categóricas incompatibles.")
417
418         if A_num.shape[1] == 0 and A_cat.shape[1] == 0:
419             return np.zeros((len(df), 0), dtype=float)
420
421     return np.hstack([A_num, A_cat])
422
423     """Celda 4 con RandomForestPortable (RF Base y RF BC_124)"""
424
425     # -----
426     # RandomForestPortable
427     # -----
428
429     class RandomForestPortable:
430         """
431             Implementación portable de un RandomForestClassifier a partir de su JSON.
432

```

```

431     - Usa PortablePreprocessor para replicar el preprocesado.
432     - Cada árbol se evalúa explícitamente sobre la matriz transformada.
433     """
434
435     def __init__(self, artifact: Dict[str, Any]) -> None:
436         """
437             Parámetros
438             -----
439             artifact : dict
440                 Diccionario cargado desde 'rf_base_portable.json' o
441                 'rf_124_portable.json'. Debe contener, al menos:
442
443                 - 'model_type': "RandomForestClassifier"
444                 - 'classes', 'n_classes'
445                 - 'feature_space':
446                     'expected_raw_features': {'num': [...], 'cat': [...]},
447                     'alias_mapping_used_in_train': {...},
448                     'dtypes_expect': {...},
449                     'preprocess': {...}
450                 }
451                 - 'forest':
452                     'n_estimators', 'max_depth', 'n_features_in', 'trees': [...]
453                 }
454             """
455             self.artifact = artifact
456             if artifact.get("model_type") != "RandomForestClassifier":
457                 raise ValueError("Solo se soporta RandomForestClassifier en este
458                                 portable.")
459
460             # Clases y metadatos
461             self.classes_: List[str] = [str(c) for c in artifact.get("classes", [])]
462             self.n_classes_: int = int(artifact.get("n_classes", len(self.classes_)))
463
464             # Espacio de características y procesador
465             fs = artifact.get("feature_space", {}) or {}
466             self.expected_raw_features: Dict[str, List[str]] = fs.get(
467                 "expected_raw_features", {}) or {}
468             self.alias_mapping: Dict[str, str] = fs.get("alias_mapping_used_in_train", {})
469             or {}
470             self.dtypes_expect: Dict[str, str] = fs.get("dtypes_expect", {}) or {}
471             prep_schema = fs.get("preprocess", {}) or {}
472             self.preprocessor = PortablePreprocessor(prep_schema)
473
474             # Bosque
475             forest = artifact.get("forest", {}) or {}
476             self.n_estimators: int = int(forest.get("n_estimators", 0))
477             self.max_depth: Optional[int] = (
478                 int(forest["max_depth"]) if forest.get("max_depth") is not None else None
479             )
480             self.n_features_in_: int = int(forest.get("n_features_in", 0))
481
482             self.trees: List[Dict[str, Any]] = list(forest.get("trees", []))
483             if not self.trees:
484                 raise ValueError("El modelo RF portable no contiene árboles.")
485             if self.n_classes_ <= 0:
486                 raise ValueError("El modelo RF portable no tiene clases definidas.")
487
488             # ----- adaptador de entrada -----
489
490     def _prepare_raw_input(self, df_raw: pd.DataFrame) -> pd.DataFrame:
491         """
492             Aplica alias de train y asegura que están todas las columnas
493             de expected_raw_features.num y expected_raw_features.cat.
494
495             - alias_mapping_used_in_train se interpreta como:
496                 {nombre_original_en_train: nombre_usado_en_modelo}
497             """
498             if not isinstance(df_raw, pd.DataFrame):
499                 raise TypeError("RandomForestPortable espera un pandas.DataFrame como
500                                 entrada.")
501
502             df = df_raw.copy()

```

```

499
500     # 1) Aliases de train: preferred -> used
501     #   alias_mapping_used_in_train viene típico como:
502     #   {"Criterio_precio_p": "C_precio_p", ...}
503     for preferred, used in self.alias_mapping.items():
504         if used in df.columns:
505             # Ya existe la columna con el nombre usado en train
506             continue
507         if preferred in df.columns:
508             # Creamos la columna "used" copiando la "preferred"
509             df[used] = df[preferred]
510
511     cols_num = list(self.expected_raw_features.get("num", []))
512     cols_cat = list(self.expected_raw_features.get("cat", []))
513     ordered = cols_num + cols_cat
514
515     # 2) Asegurar columnas y tipos mínimos
516     for c in ordered:
517         if c not in df.columns:
518             # Columna esperada que no está: la creamos como NaN
519             df[c] = np.nan
520
521     # Numéricas a float
522     for c in cols_num:
523         df[c] = pd.to_numeric(df[c], errors="coerce")
524
525     # Categóricas: no forzamos type, el PortablePreprocessor las tratará
526     return df[ordered]
527
528 # ----- evaluación de árboles -----
529
530 def _tree_predict_proba(self, tree: Dict[str, Any], X: np.ndarray) -> np.ndarray:
531     """
532     Evalúa un único árbol sobre todas las filas de X y devuelve
533     una matriz (n_samples, n_classes) con las probabilidades de ese árbol.
534     """
535     children_left = np.asarray(tree["children_left"], dtype=int)
536     children_right = np.asarray(tree["children_right"], dtype=int)
537     feature = np.asarray(tree["feature"], dtype=int)
538     threshold = np.asarray(tree["threshold"], dtype=float)
539     value = np.asarray(tree["value"], dtype=float) # shape (n_nodes, n_classes)
540
541     n_nodes = int(tree.get("n_nodes", len(children_left)))
542     n_samples = X.shape[0]
543
544     out = np.zeros((n_samples, self.n_classes_), dtype=float)
545
546     for i in range(n_samples):
547         node = 0
548         # Descenso por el árbol
549         while True:
550             # Nodo fuera de rango: por robustez, volvemos a la raíz y salimos
551             if node < 0 or node >= n_nodes:
552                 node = 0
553                 break
554
555             # Consideramos hoja si no tiene hijos o si la feature es < 0
556             if (
557                 children_left[node] == -1
558                 and children_right[node] == -1
559             ) or feature[node] < 0:
560                 break
561
562             f = feature[node]
563             thr = threshold[node]
564             x_ij = X[i, f]
565
566             if x_ij <= thr:
567                 node = children_left[node]
568             else:
569                 node = children_right[node]

```

```

571     if node < 0 or node >= n_nodes:
572         # Nodo inválido: por robustez, usamos la raíz
573         node = 0
574
575         proba_node = value[node]
576         # Re-normalizamos por seguridad
577         s = proba_node.sum()
578         if s <= 0.0:
579             proba_node = np.ones(self.n_classes_, dtype=float) / float(self.
580             n_classes_)
581         else:
582             proba_node = proba_node / s
583
584         out[i, :] = proba_node
585
586     return out
587
588     # ----- API pública -----
589
590     def predict_proba(self, df_raw: pd.DataFrame) -> np.ndarray:
591         """
592             Devuelve la matriz de probabilidades (n_samples, n_classes),
593             promediando las predicciones de todos los árboles.
594         """
595         df_prep = self._prepare_raw_input(df_raw)
596         X = self.preprocessor.transform(df_prep)
597         X = _ensure_2d_array(X)
598
599         if X.shape[1] != self.n_features_in_ and self.n_features_in_ > 0:
600             # Aviso suave; no lanzamos error duro para permitir cierta tolerancia
601             pass
602
603         n_samples = X.shape[0]
604         proba_acc = np.zeros((n_samples, self.n_classes_), dtype=float)
605
606         for tree in self.trees:
607             proba_acc += self._tree_predict_proba(tree, X)
608
609         # Promedio de todos los árboles
610         proba_acc /= float(len(self.trees))
611
612         # Normalización final por seguridad
613         sums = proba_acc.sum(axis=1, keepdims=True)
614         sums[sums == 0.0] = 1.0
615         proba_acc /= sums
616
617     return proba_acc
618
619     def predict(self, df_raw: pd.DataFrame) -> np.ndarray:
620         """
621             Devuelve las clases predichas (np.ndarray de str).
622         """
623         proba = self.predict_proba(df_raw)
624         idx = np.argmax(proba, axis=1)
625         classes = np.array(self.classes_, dtype=object)
626         return classes[idx]
627
628     """Celda 5 - XGBBajaPortable (XGBoost multi-clase, umbral BAJO fijo 0,38)"""
629
630     # -----
631     # XGBBajaPortable
632     # -----
633
634     class XGBBajaPortable:
635         """
636             Implementación portable del modelo XGBoost multi-clase (BAJO/MEDIO/ALTO)
637             para predecir baja_comp.
638
639             - Objetivo: multi:softprob.
640             - A partir de xgb_baja_portable.json + xgb_baja_model.json.
641         """
642
643         def __init__(self, meta: Dict[str, Any], booster: Dict[str, Any]) -> None:
644             self.meta = meta

```

```

642     self.variables_predictoras: List[str] = list(meta.get("variables_predictoras",
643         []))
644     self.encoders_schema: Dict[str, Dict[str, Any]] = meta.get("encoders_schema",
645         {}) or {}
646     self.labels: List[str] = list(meta.get("labels", [])) or ["BAJO", "MEDIO",
647         "ALTO"]
648     self.policies: Dict[str, Any] = meta.get("policies", {}) or {}
649     self.fallback_class_por_col: Dict[str, str] = self.policies.get(
650         "fallback_class_por_col", {})
651
652     # Umbral fijo para la clase BAJO, independientemente del meta
653     # (decisión explícita para este script portable).
654     self.threshold_bajo: float = 0.38
655
656     # Booster JSON (estructura interna de XGBoost)
657     learner = booster.get("learner", {}) or {}
658     self.feature_names: List[str] = list(learner.get("feature_names", []))
659     lm_param = learner.get("learner_model_param", {}) or {}
660     self.base_score: float = float(lm_param.get("base_score", 0.5))
661     self.num_class: int = int(lm_param.get("num_class", len(self.labels) or 3))
662
663     gb = learner.get("gradient_booster", {}) or {}
664     model = gb.get("model", {}) or {}
665     self.trees: List[Dict[str, Any]] = list(model.get("trees", []))
666     self.tree_info: List[int] = list(model.get("tree_info", []))
667
668     if len(self.tree_info) != len(self.trees):
669         raise ValueError("Inconsistencia entre número de árboles y tree_info en
670                           XGB.")
671
672     # Índice de la clase BAJO
673     try:
674         self.idx_bajo: int = self.labels.index("BAJO")
675     except ValueError:
676         self.idx_bajo = 0
677
678     # Preconstruimos esquemas de encoders: col -> {clase: índice}
679     self._encoders_index: Dict[str, Dict[str, int]] = {}
680     for col, schema in self.encoders_schema.items():
681         classes = [str(c) for c in schema.get("classes", [])]
682         self._encoders_index[col] = {cls: i for i, cls in enumerate(classes)}
683
684     # ----- preprocessado de entrada -----
685
686     def _prepare_X(self, df_raw: pd.DataFrame) -> pd.DataFrame:
687         """
688             Construye el DataFrame X con las variables_predictoras en el orden
689             exacto esperado por el booster.
690         """
691
692         X = pd.DataFrame(index=df_raw.index)
693
694         # 1) Columnas básicas (num y cat ya codificados como en df_norm)
695         for col in self.variables_predictoras:
696             if col in df_raw.columns:
697                 X[col] = df_raw[col]
698             else:
699                 # Política de columna numérica faltante: fill_0 (coherente con meta)
700                 X[col] = 0
701
702         # 2) Encoders categóricos tipo LabelEncoder
703         fallback_map = self.fallback_class_por_col or {}
704         for col, idx_map in self._encoders_index.items():
705             if col in X.columns:
706                 s = X[col].astype(str)
707                 classes = set(idx_map.keys())
708                 fb = str(fallback_map.get(col, next(iter(classes)) if classes else ""))
709

```

```

708     # Sustituimos valores fuera del repertorio por el fallback
709     s = s.where(s.isin(classes), other=fb)
710
711     # Transformamos a índices
712     def to_idx(val: Any) -> int:
713         try:
714             return idx_map[str(val)]
715         except KeyError:
716             return idx_map.get(fb, 0)
717
718         X[col] = s.map(to_idx).astype(int)
719     else:
720         # Columna prevista pero ausente -> 0
721         X[col] = 0
722
723
724     # 3) Mes_lici a entero si está
725     if "Mes_lici" in X.columns and not pd.api.types.is_integer_dtype(X["Mes_lici"]):
726         X["Mes_lici"] = (
727             pd.to_numeric(X["Mes_lici"], errors="coerce")
728             .fillna(0)
729             .astype(int)
730         )
731
732     return X[self.variables_predictoras]
733
734 # ----- evaluación de árboles -----
735
736 def _tree_predict_margin(self, tree: Dict[str, Any], X: np.ndarray) -> np.ndarray:
737 """
738     Evalúa un árbol XGBoost sobre X y devuelve los márgenes (una dimensión).
739 """
740     left_children = np.asarray(tree["left_children"], dtype=int)
741     right_children = np.asarray(tree["right_children"], dtype=int)
742     split_indices = np.asarray(tree["split_indices"], dtype=int)
743     split_conditions = np.asarray(tree["split_conditions"], dtype=float)
744     default_left = np.asarray(tree["default_left"], dtype=int)
745     base_weights = np.asarray(tree["base_weights"], dtype=float)
746
747     n_nodes = int(tree.get("tree_param", {}).get("num_nodes", len(left_children)))
748     n_samples = X.shape[0]
749
750     out = np.zeros(n_samples, dtype=float)
751
752     for i in range(n_samples):
753         node = 0
754         while node >= 0 and node < n_nodes:
755             left = left_children[node]
756             right = right_children[node]
757
758             # Hoja: no hay hijos
759             if left == -1 and right == -1:
760                 break
761
762             feat_idx = split_indices[node]
763             thr = split_conditions[node]
764             v = X[i, feat_idx]
765
766             if math.isnan(v):
767                 # Usamos default_left
768                 go_left = bool(default_left[node])
769             else:
770                 # Regla típica de XGBoost: v < thr -> izquierda
771                 go_left = v < thr
772
773             node = left if go_left else right
774             if node == -1:
775                 break
776
777             if node == -1:
778                 # Por robustez, volvemos a la raíz

```

```

779         node = 0
780
781         out[i] = base_weights[node]
782
783     return out
784
785 # ----- API pública -----
786
787 def predict_proba(self, df_raw: pd.DataFrame) -> np.ndarray:
788     """
789     Devuelve la matriz de probabilidades (n_samples, num_class)
790     del XGBoost multi-clase (BAJO/MEDIO/ALTO).
791     """
792     X_df = self._prepare_X(df_raw)
793     X = X_df.to_numpy(dtype=float)
794     X = _ensure_2d_array(X)
795
796     n_samples = X.shape[0]
797     margins = np.full((n_samples, self.num_class), self.base_score, dtype=float)
798
799     for t_idx, tree in enumerate(self.trees):
800         class_id = int(self.tree_info[t_idx])
801         if class_id < 0 or class_id >= self.num_class:
802             continue
803         contrib = self._tree_predict_margin(tree, X)
804         margins[:, class_id] += contrib
805
806     # Usamos el softmax fila a fila definido en las utilidades
807     proba = _softmax(margins)
808     return proba
809
810 def predict_baja(self, df_raw: pd.DataFrame) -> Tuple[pd.Series, pd.Series]:
811     """
812     Devuelve:
813         - baja_comp: Series de int (0/1)
814         - p_bajo: Series de float (probabilidad de BAJO)
815     """
816     proba = self.predict_proba(df_raw)
817     if proba.shape[1] <= self.idx_bajo:
818         # Modelo mal formado; devolvemos todo ceros
819         p_bajo = np.zeros(proba.shape[0], dtype=float)
820     else:
821         p_bajo = proba[:, self.idx_bajo]
822
823     baja_flag = (p_bajo >= self.threshold_bajo).astype(np.int8)
824     idx = df_raw.index
825     return (
826         pd.Series(baja_flag, index=idx, name="baja_comp"),
827         pd.Series(p_bajo, index=idx, name="p_bajo"),
828     )
829
830 """Celda 6 - KNNPortable (espacio PRE_NUM, distancias y similitud)"""
831
832 # -----
833 # KNNPortable
834 # -----
835
836 class KNNPortable:
837     """
838         Implementación portable del espacio kNN PRE_NUM.
839
840         Se basa en:
841             - knn_space_portable.json
842             - knn_train.npz (Z_train)
843             - knn_all.npz    (Z_all)
844
845         Funcionalidad:
846             - Proyectar nuevas filas al espacio PRE_NUM con el pipeline numérico
847                 (imputación + escalado).
848             - Calcular distancias euclídeas a Z_train o Z_all.
849             - Devolver índices y distancias de los k vecinos, junto con una
850                 medida de similitud (%) basada en p95_nn1_train.

```

```

851 """
852
853     def __init__(self,
854         meta: Dict[str, Any],
855         Z_train: np.ndarray,
856         Z_all: np.ndarray,
857     ) -> None:
858         self.meta = meta
859         self.pre_num_features: List[str] = list(meta.get("pre_num_features", []))
860         self.metric: str = str(meta.get("metric", "euclidean"))
861         self.k_neighbors: int = int(meta.get("k_neighbors", 25))
862         self.p95_nn1_train: float = float(meta.get("p95_nn1_train", 0.0))
863
864         num_pipeline = meta.get("num_pipeline", {}) or {}
865         num_imp = num_pipeline.get("imputer", {}) or {}
866         num_scl = num_pipeline.get("scaler", {}) or {}
867
868         self.imp_strategy: str = num_imp.get("strategy", "median")
869         self.imp_statistics: np.ndarray = np.asarray(
870             num_imp.get("statistics", []), dtype=float
871         )
872
873         self.with_mean: bool = bool(num_scl.get("with_mean", True))
874         self.with_std: bool = bool(num_scl.get("with_std", True))
875         self.mean_: np.ndarray = np.asarray(
876             num_scl.get("mean", np.zeros(len(self.pre_num_features))), dtype=float
877         )
878
879         self.scale_: np.ndarray = np.asarray(
880             num_scl.get("scale", np.ones(len(self.pre_num_features))), dtype=float
881         )
882         self.scale_[self.scale_ == 0.0] = 1.0 # evitar divisiones por cero
883
884         # Catálogo (metadatos informativos; no podemos reconstruir el DF desde aquí)
885         self.catalog_meta: Dict[str, Any] = meta.get("catalog", {}) or {}
886
887         # Matrices de referencia
888         self.Z_train: np.ndarray = np.asarray(Z_train, dtype=float)
889         self.Z_all: np.ndarray = np.asarray(Z_all, dtype=float)
890
891         if self.Z_train.ndim != 2 or self.Z_all.ndim != 2:
892             raise ValueError("Z_train y Z_all deben ser matrices 2D.")
893
894 # ----- proyección al espacio PRE_NUM -----
895
896     def build_X_num(self, df: pd.DataFrame) -> pd.DataFrame:
897         """
898             Construye el DataFrame numérico X_num con las columnas pre_num_features
899             en el orden exacto, aplicando algunas reglas de fallback coherentes
900             con el pipeline original (C_precio_p <- Criterio_precio_p).
901         """
902         if not isinstance(df, pd.DataFrame):
903             raise TypeError("KNNPortable espera un pandas.DataFrame como entrada.")
904
905         X_num = pd.DataFrame(index=df.index)
906
907         for c in self.pre_num_features:
908             if c in df.columns:
909                 X_num[c] = pd.to_numeric(df[c], errors="coerce")
910             elif c == "C_precio_p" and "Criterio_precio_p" in df.columns:
911                 # Fallback: usar la columna original de criterio si existe
912                 X_num[c] = pd.to_numeric(df["Criterio_precio_p"], errors="coerce")
913             else:
914                 X_num[c] = np.nan
915
916         return X_num[self.pre_num_features]
917
918     def _transform(self, df: pd.DataFrame) -> np.ndarray:
919         """
920             Aplica el imputer + scaler numérico del pipeline kNN a df, y devuelve
921             la matriz NumPy 2D en el espacio PRE_NUM.
922         """

```

```

923     X_num = self._build_X_num(df)
924     A = X_num.to_numpy(dtype=float)
925
926     # Imputación por estadísticos
927     if self.imp_strategy in ("median", "most_frequent"):
928         stats = self.imp_statistics
929         if stats.shape[0] != A.shape[1]:
930             raise ValueError(
931                 "Dimensión de estadísticas de imputación kNN inconsistente: "
932                 f"{stats.shape[0]} != {A.shape[1]}")
933     )
934     for j in range(A.shape[1]):
935         col = A[:, j]
936         mask = np.isnan(col)
937         if mask.any():
938             col[mask] = stats[j]
939
940     # Escalado
941     if self.with_mean:
942         A = A - self.mean_
943     if self.with_std:
944         A = A / self.scale_
945
946     return A
947
948 # ----- distancias y vecinos -----
949
950 def _compute_distances(self, Zq: np.ndarray, space: str) -> np.ndarray:
951     """
952     Calcula la matriz de distancias entre Zq (consultas) y Z_base
953     (train o all) usando la métrica configurada.
954
955     Actualmente solo se implementa 'euclidean'; otras métricas se
956     degradan a euclidiana de forma explícita.
957     """
958     if space == "train":
959         Z_base = self.Z_train
960     elif space == "all":
961         Z_base = self.Z_all
962     else:
963         raise ValueError("space debe ser 'train' o 'all'.")
964
965     if self.metric not in ("euclidean", "l2"):
966         # Degradamos a euclidiana, lo documentamos en el docstring
967         pass
968
969     # Broadcasting: (n_q, 1, d) - (1, n_ref, d) -> (n_q, n_ref, d)
970     diff = Zq[:, None, :] - Z_base[None, :, :]
971     dist = np.sqrt(np.sum(diff ** 2, axis=2)) # shape (n_q, n_ref)
972     return dist
973
974 def neighbors(
975     self,
976     df: pd.DataFrame,
977     topk: int = 10,
978     space: str = "all",
979 ) -> Dict[str, Any]:
980     """
981     Calcula los k vecinos más cercanos para cada fila de df.
982
983     Parámetros
984     -----
985     df : DataFrame
986         DataFrame con las variables PRE necesarias (al menos las que
987         figuran en pre_num_features).
988     topk : int, por defecto 10
989         Número de vecinos a devolver (se trunca si hay menos filas
990         en Z_base).
991     space : {"train", "all"}, por defecto "all"
992         Conjunto de referencia (Z_train o Z_all).
993
994     Devuelve

```

```

995     -----
996     dict con:
997         - "indices" : np.ndarray (n_samples, k) de índices de vecinos.
998         - "distances": np.ndarray (n_samples, k) de distancias euclídeas.
999         - "similarity_percent": np.ndarray (n_samples, k) con sim% = max(0, 1 -
d/p95)*100.
1000         - "p95_nn1_train": float (el mismo valor meta).
1001     """
1002     Zq = self._transform(df)
1003     Zq = _ensure_2d_array(Zq)
1004
1005     dist = self._compute_distances(Zq, space=space)
1006     n_q, n_ref = dist.shape
1007
1008     k = max(1, min(int(topk), n_ref))
1009     # argsort a lo largo del eje de referencia
1010     # índices de vecinos ordenados por distancia ascendente
1011     idx_sorted = np.argsort(dist, axis=1)[:, :k]
1012
1013     # Distancias ordenadas
1014     rows = np.arange(n_q)[:, None]
1015     dist_sorted = dist[rows, idx_sorted]
1016
1017     # Similitud porcentual basada en p95_nn1_train
1018     if self.p95_nn1_train > 0:
1019         sim_percent = np.maximum(
1020             0.0,
1021             1.0 - dist_sorted / float(self.p95_nn1_train),
1022         ) * 100.0
1023     else:
1024         sim_percent = np.zeros_like(dist_sorted)
1025
1026     return {
1027         "indices": idx_sorted.astype(int),
1028         "distances": dist_sorted.astype(float),
1029         "similarity_percent": sim_percent.astype(float),
1030         "p95_nn1_train": float(self.p95_nn1_train),
1031     }
1032
1033 """Celda 7 - EnsemblePredictor (nueva lógica RF base vs RF_124)"""
1034
1035 # -----
1036 # EnsemblePredictor
1037 # -----
1038
1039 class EnsemblePredictor:
1040     """
1041     Ensemble simplificado entre RF base y RF_124_BASE+BC.
1042
1043     Lógica:
1044         - Tomamos RF base como predicción por defecto.
1045         - Solo activamos el ensemble cuando la clase top-1 del RF base está
1046           en classes_124 = {"1", "2", "4"}.
1047         - Además exigimos que la clase top-1 del RF_124 también esté en ese
1048           mismo conjunto.
1049         - En ese caso, comparamos la probabilidad top-1 (p1) de cada modelo
1050           y escogemos el modelo cuya p1 sea mayor.
1051         - En caso de empate, se mantiene la predicción del RF base.
1052         - margin_min y otras reglas "seguras" NO se aplican aquí.
1053
1054     La salida incorpora también la información de baja_comp del modelo XGB.
1055     """
1056
1057     def __init__(self, classes_124: Optional[List[str]] = None) -> None:
1058         if classes_124 is None:
1059             classes_124 = ["1", "2", "4"]
1060         # Guardamos el conjunto de clases donde se considera el especialista
1061         self.classes_124 = set(str(c) for c in classes_124)
1062
1063     # ----- utilidades internas -----
1064
1065     @staticmethod

```

```

1066
1067     def _top1_top2(row: np.ndarray, labels: List[str]) -> Tuple[str, float, float]:
1068         """
1069             Dada una fila de probabilidades y la lista de labels,
1070             devuelve (label_top1, p1, p2), donde p2 es la segunda probabilidad
1071             más alta (útil si quisieras volver a un criterio por margen).
1072         """
1073         row = np.asarray(row, dtype=float)
1074         if row.ndim != 1:
1075             raise ValueError("_top1_top2 espera un vector 1D de probabilidades.")
1076
1077         # Ordenamos índices por probabilidad descendente
1078         idx_sorted = np.argsort(row)[::-1]
1079         idx1 = idx_sorted[0]
1080         p1 = float(row[idx1])
1081
1082         if len(idx_sorted) > 1:
1083             idx2 = idx_sorted[1]
1084             p2 = float(row[idx2])
1085         else:
1086             p2 = 0.0
1087
1088         label_top1 = str(labels[idx1])
1089         return label_top1, p1, p2
1090
1091     # ----- API principal -----
1092
1093     def apply(
1094         self,
1095         pred_base: Any,
1096         proba_base: Any,
1097         labels_base: List[str],
1098         pred_124: Any,
1099         proba_124: Any,
1100         labels_124: List[str],
1101         baja_comp: pd.Series,
1102         p_bajo: pd.Series,
1103     ) -> pd.DataFrame:
1104         """
1105             Aplica la lógica de ensemble sobre las predicciones de RF base,
1106             RF_124 y XGB baja.
1107
1108             Parámetros
1109             -----
1110             pred_base : array-like
1111                 Clases predichas por RF base (no se usan directamente aquí,
1112                 se recalculan desde las probabilidades).
1113             proba_base : array-like (n_samples, n_classes_base)
1114                 Probabilidades de RF base.
1115             labels_base : list[str]
1116                 Etiquetas de las clases de RF base en el mismo orden que las columnas
1117                 de proba_base.
1118             pred_124 : array-like
1119                 Clases predichas por RF_124 (no se usan directamente aquí).
1120             proba_124 : array-like (n_samples, n_classes_124)
1121                 Probabilidades de RF_124.
1122             labels_124 : list[str]
1123                 Etiquetas de las clases de RF_124 en el mismo orden que las columnas
1124                 de proba_124.
1125             baja_comp : pd.Series
1126                 Índicador 0/1 de baja_comp procedente de XGB.
1127             p_bajo : pd.Series
1128                 Probabilidad de clase BAJO procedente de XGB.
1129
1130             Devuelve
1131             -----
1132             DataFrame con columnas:
1133                 - 'pred'           : clase final (str)
1134                 - 'confianza_pred' : probabilidad top-1 del modelo elegido
1135                 - 'fuente_confianza' : 'rf_base' o 'rf_124_basebc'
1136                 - 'flip_ensemble'   : bool, True si se cambió desde rf_base a rf_124
1137                 - 'baja_comp'       : 0/1
1138                 - 'p_bajo'          : probabilidad de BAJO

```

```

1138 """
1139 # Convertimos a arrays NumPy
1140 proba_base = np.asarray(proba_base, dtype=float)
1141 proba_124 = np.asarray(proba_124, dtype=float)
1142
1143 n = proba_base.shape[0]
1144 if proba_124.shape[0] != n:
1145     raise ValueError("proba_base y proba_124 tienen número de filas distinto.")
1146
1147 # Aseguramos que baja_comp y p_bajo tienen longitud compatible
1148 if len(baja_comp) != n or len(p_bajo) != n:
1149     raise ValueError(
1150         "baja_comp y p_bajo deben tener la misma longitud que las matrices de
1151         probas."
1152     )
1153
1154 # Index para el DataFrame de salida
1155 idx_out = baja_comp.index
1156
1157 # Inicializamos salidas
1158 pred_final = np.empty(n, dtype=object)
1159 fuente = np.full(n, "rf_base", dtype=object)
1160 conf_final = np.zeros(n, dtype=float)
1161 flip = np.zeros(n, dtype=bool)
1162
1163 # Recorremos fila a fila
1164 for i in range(n):
1165     row_base = proba_base[i]
1166     row_124 = proba_124[i]
1167
1168     base_label_top1, p1_base, _ = self._top1_top2(row_base, labels_base)
1169     spec_label_top1, p1_spec, _ = self._top1_top2(row_124, labels_124)
1170
1171     # Por defecto, siempre usamos RF base
1172     pred_final[i] = base_label_top1
1173     conf_final[i] = p1_base
1174     fuente[i] = "rf_base"
1175
1176     # Solo activamos ensemble si la clase base está en {1,2,4}
1177     if base_label_top1 not in self.classes_124:
1178         continue
1179
1180     # Y si RF_124 también predice una clase válida del especialista
1181     if spec_label_top1 not in self.classes_124:
1182         continue
1183
1184     # Regla de decisión: elegir el modelo con mayor probabilidad top-1
1185     if p1_spec > p1_base:
1186         pred_final[i] = spec_label_top1
1187         conf_final[i] = p1_spec
1188         fuente[i] = "rf_124_basebc"
1189         flip[i] = True
1190
1191     # Si p1_spec <= p1_base, nos quedamos con la predicción base
1192
1193 df_out = pd.DataFrame(
1194     {
1195         "pred": pred_final.astype(str),
1196         "confianza_pred": conf_final.astype(float),
1197         "fuente_confianza": fuente.astype(str),
1198         "flip_ensemble": flip.astype(bool),
1199         "baja_comp": baja_comp.astype(int).values,
1200         "p_bajo": p_bajo.astype(float).values,
1201     },
1202     index=idx_out,
1203 )
1204
1205 return df_out
1206
1207 """
1208 """Celda 8 - PortablePredictor + load_predictor_from_zip (API de alto nivel)"""
1209

```

```
1208 # PortablePredictor: fachada principal
1209 # -----
1210
1211 class PortablePredictor:
1212     """
1213         Fachada principal para usar todos los modelos portables desde un ZIP.
1214
1215     Flujo típico de uso:
1216         predictor = PortablePredictor.from_zip("modelos_portable.zip")
1217         resultados = predictor.predict(df_norm)
1218         vecinos = predictor.knn_neighbors(df_norm.iloc[[0]], topk=10, space="all")
1219
1220     Importante:
1221         - Este predictor asume que df_norm (df_raw en las firmas) ya está en el
1222             formato "normalizado" que consumen los modelos:
1223             * columnas y tipos compatibles con expected_raw_features de los RF
1224             * variables_predictoras del XGB de baja_comp
1225             * pre_num_features del espacio kNN
1226             Es decir, df_norm debería ser la salida del script de transformación
1227             (export_python_transformacion.py) aplicada sobre los datos originales.
1228     """
1229
1230     def __init__(
1231         self,
1232         zip_path: str,
1233         store: ZipStore,
1234         rf_base: RandomForestPortable,
1235         rf_124: RandomForestPortable,
1236         xgb_baja: XGBBajaPortable,
1237         knn: KNNPortable,
1238         ensemble: EnsemblePredictor,
1239         baja_comp_dup_times: int,
1240         ensemble_meta: Dict[str, Any],
1241         preprocess_pre: Optional[PortablePreprocessor] = None,
1242         centroids_pre: Optional[Dict[str, Any]] = None,
1243     ) -> None:
1244         self.zip_path = zip_path
1245         self.store = store
1246         self.rf_base = rf_base
1247         self.rf_124 = rf_124
1248         self.xgb_baja = xgb_baja
1249         self.knn = knn
1250         self.ensemble = ensemble
1251         self.baja_comp_dup_times = int(baja_comp_dup_times)
1252         self.ensemble_meta = ensemble_meta
1253         self.preprocess_pre = preprocess_pre
1254         self.centroids_pre = centroids_pre
1255
1256     # ----- construcción desde ZIP -----
1257
1258     @classmethod
1259     def from_zip(cls, zip_path: str) -> "PortablePredictor":
1260         """
1261             Carga todos los artefactos portables desde un ZIP y construye
1262             un PortablePredictor listo para usar.
1263
1264             Parámetros
1265             -----
1266             zip_path : str
1267                 Ruta al fichero modelos_portable.zip.
1268
1269             Devuelve
1270             -----
1271             PortablePredictor
1272         """
1273         store = ZipStore(zip_path)
1274
1275         # 1) Meta del ensemble (punto de entrada común)
1276         if not store.has("ensemble_portable.json"):
1277             raise FileNotFoundError(
1278                 "No se ha encontrado 'ensemble_portable.json' en el ZIP."
1279         )
```

```

1280 ensemble_meta = store.read_json("ensemble_portable.json")
1281
1282 deps = ensemble_meta.get("depends_on_portable", {}) or {}
1283
1284 # 2) Artefactos RF base y RF_124
1285 rf_base_name = deps.get("rf_base", "rf_base_portable.json")
1286 rf_124_name = deps.get("rf_124_basebc", "rf_124_portable.json")
1287
1288 if not store.has(rf_base_name):
1289     raise FileNotFoundError(f"Fallta artefacto RF base: {rf_base_name!r}")
1290 if not store.has(rf_124_name):
1291     raise FileNotFoundError(f"Fallta artefacto RF_124: {rf_124_name!r}")
1292
1293 rf_base_art = store.read_json(rf_base_name)
1294 rf_124_art = store.read_json(rf_124_name)
1295
1296 rf_base = RandomForestPortable(rf_base_art)
1297 rf_124 = RandomForestPortable(rf_124_art)
1298
1299 # 3) Artefacto XGB baja_comp (meta + booster)
1300 xgb_meta_name = deps.get("xgb_baja", "xgb_baja_portable.json")
1301 if not store.has(xgb_meta_name):
1302     raise FileNotFoundError(f"Fallta artefacto XGB meta: {xgb_meta_name!r}")
1303 xgb_meta = store.read_json(xgb_meta_name)
1304
1305 # Booster guardado con Booster.save_model en formato JSON
1306 if not store.has("xgb_baja_model.json"):
1307     raise FileNotFoundError("Fallta 'xgb_baja_model.json' en el ZIP.")
1308 xgb_booster = store.read_json("xgb_baja_model.json")
1309
1310 xgb_baja = XGBBajaPortable(xgb_meta, xgb_booster)
1311
1312 # 4) Espacio kNN PRE_NUM
1313 knn_meta_name = deps.get("knn_space", "knn_space_portable.json")
1314 if not store.has(knn_meta_name):
1315     raise FileNotFoundError(f"Fallta artefacto kNN meta: {knn_meta_name!r}")
1316 knn_meta = store.read_json(knn_meta_name)
1317
1318 if not store.has("knn_train.npz"):
1319     raise FileNotFoundError("Fallta 'knn_train.npz' en el ZIP.")
1320 if not store.has("knn_all.npz"):
1321     raise FileNotFoundError("Fallta 'knn_all.npz' en el ZIP.")
1322
1323 Z_train_npz = store.read_npz("knn_train.npz")
1324 Z_all_npz = store.read_npz("knn_all.npz")
1325
1326 if "Z" not in Z_train_npz.files:
1327     raise KeyError("En 'knn_train.npz' se esperaba un array 'Z'.")
1328 if "Z" not in Z_all_npz.files:
1329     raise KeyError("En 'knn_all.npz' se esperaba un array 'Z'.")
1330
1331 knn = KNNPortable(
1332     meta=knn_meta,
1333     Z_train=Z_train_npz["Z"],
1334     Z_all=Z_all_npz["Z"],
1335 )
1336
1337 # 5) Preprocesador PRE global (opcional, por si quieres usarlo fuera)
1338 preprocess_pre = None
1339 centroids_pre = None
1340 pre_pre_name = deps.get("preprocess_pre", "preprocess_pre_portable.json")
1341 if store.has(pre_pre_name):
1342     pre_pre_art = store.read_json(pre_pre_name)
1343     # El esquema de preprocess_pre es compatible con PortablePreprocessor
1344     prep_schema = {
1345         "numeric": pre_pre_art.get("numeric", {}),
1346         "categorical": pre_pre_art.get("categorical", {}),
1347         "feature_names_out": pre_pre_art.get("feature_names_out"),
1348     }
1349     preprocess_pre = PortablePreprocessor(prep_schema)
1350     centroids_pre = pre_pre_art.get("centroids_ref") or pre_pre_art.get(
1351         "centroids")

```

```

1351
1352     # 6) Ensemble
1353     ensemble = EnsemblePredictor(classes_124=["1", "2", "4"])
1354
1355     # 7) Configuración de baja_comp duplicada para RF_124
1356     baja_cfg = ensemble_meta.get("baja_comp", {}) or {}
1357     baja_comp_dup_times = int(baja_cfg.get("duplicated_times", 5))
1358
1359     return cls(
1360         zip_path=zip_path,
1361         store=store,
1362         rf_base=rf_base,
1363         rf_124=rf_124,
1364         xgb_baja=xgb_baja,
1365         knn=knn,
1366         ensemble=ensemble,
1367         baja_comp_dup_times=baja_comp_dup_times,
1368         ensemble_meta=ensemble_meta,
1369         preprocess_pre=preprocess_pre,
1370         centroids_pre=centroids_pre,
1371     )
1372
1373     # ----- utilidades internas -----
1374
1375     @staticmethod
1376     def _ensure_dataframe(df: Any) -> pd.DataFrame:
1377         """
1378             Asegura que la entrada se convierte en un DataFrame.
1379
1380             - Si ya es DataFrame, se devuelve una copia.
1381             - Si es dict, se interpreta como un único registro (index = [0]).
1382             - Si es Series, se interpreta como una fila (index = [0]).
1383         """
1384         if isinstance(df, pd.DataFrame):
1385             return df.copy()
1386         elif isinstance(df, pd.Series):
1387             return df.to_frame().T.copy()
1388         elif isinstance(df, dict):
1389             return pd.DataFrame([df]).copy()
1390         else:
1391             raise TypeError(
1392                 "PortablePredictor espera un DataFrame, Series o dict como entrada."
1393             )
1394
1395     def _attach_baja_comp_dups_for_rf124(
1396         self,
1397         df: pd.DataFrame,
1398         baja_comp: pd.Series,
1399     ) -> pd.DataFrame:
1400         """
1401             Inserta 'baja_comp' y sus duplicadas en df para el RF_124, siguiendo
1402             la información de expected_raw_features del modelo portable.
1403
1404             Estrategia:
1405                 - Columna base: 'baja_comp' (si no existe, se crea).
1406                 - Duplicadas: todas las columnas categóricas de RF_124 cuyo nombre
1407                     empiece por 'baja_comp_dup'.
1408         """
1409         df2 = df.copy()
1410
1411         # Base
1412         df2["baja_comp"] = baja_comp.astype(int)
1413
1414         # Duplicadas según expected_raw_features del RF_124
1415         exp_124 = self.rf_124.artifact.get("feature_space", {}).get(
1416             "expected_raw_features", {} or {}
1417         )
1418         cat_124 = list(exp_124.get("cat", []))
1419         dup_names = [c for c in cat_124 if c.startswith("baja_comp_dup")]
1420
1421         for dn in dup_names:
1422             df2[dn] = df2["baja_comp"]

```

```

1422     return df2
1423
1424 # ----- API pública: predicciones -----
1425
1426 def predict(self, df_raw: Any) -> Dict[str, Any]:
1427     """
1428     Ejecuta toda la cadena de predicción:
1429
1430     - XGB baja_comp -> baja_comp (0/1) + p_bajo
1431     - RF base       -> clases + probabilidades
1432     - RF_124         -> clases + probabilidades (usando baja_comp y duplicadas)
1433     - Ensemble       -> combinación RF base vs RF_124 + baja_comp
1434
1435     Parámetros
1436     -----
1437     df_raw : DataFrame, Series o dict
1438         df_norm con todas las columnas necesarias para los modelos.
1439
1440     Devuelve
1441     -----
1442     dict con:
1443     - "rf_base" : {"pred": Series, "proba": DataFrame}
1444     - "rf_124" : {"pred": Series, "proba": DataFrame}
1445     - "xgb_baja": {"baja_comp": Series, "p_bajo": Series, "threshold_bajo": float}
1446     - "ensemble": DataFrame con la predicción final y metadatos.
1447 """
1448 df = self._ensure_dataframe(df_raw)
1449
1450 # 1) XGB baja_comp
1451 baja_comp, p_bajo = self.xgb_baja.predict_baja(df)
1452
1453 # 2) RF base
1454 proba_base = self.rf_base.predict_proba(df)
1455 pred_base = self.rf_base.predict(df)
1456
1457 # 3) RF_124 (requiere baja_comp + duplicadas)
1458 df_rf124 = self._attach_baja_comp_dups_for_rf124(df, baja_comp)
1459 proba_124 = self.rf_124.predict_proba(df_rf124)
1460 pred_124 = self.rf_124.predict(df_rf124)
1461
1462 # 4) Ensemble
1463 ensemble_df = self.ensemble.apply(
1464     pred_base=pred_base,
1465     proba_base=proba_base,
1466     labels_base=self.rf_base.classes_,
1467     pred_124=pred_124,
1468     proba_124=proba_124,
1469     labels_124=self.rf_124.classes_,
1470     baja_comp=baja_comp,
1471     p_bajo=p_bajo,
1472 )
1473
1474 # Construimos salidas en estructuras cómodas
1475 idx = df.index
1476 rf_base_out = {
1477     "pred": pd.Series(pred_base, index=idx, name="rf_base_pred"),
1478     "proba": pd.DataFrame(proba_base, index=idx, columns=self.rf_base.classes_),
1479 }
1480 rf_124_out = {
1481     "pred": pd.Series(pred_124, index=idx, name="rf_124_pred"),
1482     "proba": pd.DataFrame(proba_124, index=idx, columns=self.rf_124.classes_),
1483 }
1484 xgb_out = {
1485     "baja_comp": baja_comp,
1486     "p_bajo": p_bajo,
1487     "threshold_bajo": float(self.xgb_baja.threshold_bajo),
1488 }
1489
1490 return {
1491     "rf_base": rf_base_out,

```

```

1492     "rf_124": rf_124_out,
1493     "xgb_baja": xgb_out,
1494     "ensemble": ensemble_df,
1495   }
1496
1497   def knn_neighbors(
1498     self,
1499     df: Any,
1500     topk: int = 10,
1501     space: str = "all",
1502   ) -> Dict[str, Any]:
1503     """
1504       Calcula vecinos más cercanos en el espacio PRE_NUM a partir de df.
1505
1506       Parámetros
1507       -----
1508       df : DataFrame, Series o dict
1509         df_norm con, al menos, las columnas pre_num_features del kNN.
1510       topk : int
1511         Número máximo de vecinos a devolver.
1512       space : {"train", "all"}
1513         Conjunto de referencia (Z_train o Z_all).
1514
1515       Devuelve
1516       -----
1517       dict como el de KNNPortable.neighbors().
1518       """
1519       df2 = self._ensure_dataframe(df)
1520       return self.knn.neighbors(df2, topk=topk, space=space)
1521
1522
1523 # -----#
1524 # Función de conveniencia
1525 # -----#
1526
1527 def load_predictor_from_zip(zip_path: str) -> PortablePredictor:
1528   """
1529     Punto de entrada de alto nivel para construir un PortablePredictor
1530     a partir de un ZIP de artefactos portables.
1531   """
1532   return PortablePredictor.from_zip(zip_path)
1533
1534
1535 # -----#
1536 # Ejemplo de uso (comentado)
1537 # -----#
1538
1539 # Ejemplo de cómo usarlo en tu entorno (descomentar y adaptar rutas):
1540 #
1541 # zip_path =
1542 #   "/content/drive/MyDrive/_Pipeline_produccion_prediccion/modelos_portable.zip"
1543 # predictor = load_predictor_from_zip(zip_path)
1544 #
1545 # import pandas as pd
1546 # df_test = pd.read_excel(
1547 #   "/content/drive/MyDrive/_Pipeline_produccion/prediccion/datos/datos_v8.xlsx"
1548 # )
1549 # # Supongamos que df_test ya se ha transformado con export_python_transformacion.py
1550 # # para obtener df_norm; aquí usarías df_norm en lugar de df_test si lo tienes:
1551 # df_norm = df_test # sustitúyelo por tu df_norm real
1552 #
1553 # resultados = predictor.predict(df_norm.head(5))
1554 # print(resultados["ensemble"].head())
1555 # vecinos = predictor.knn_neighbors(df_norm.iloc[[0]], topk=10, space="all")
1556 # print(vecinos["indices"][0], vecinos["distances"][0])
1557 """Celda 9 - Helper build_prediction_payload + ejemplo JSON"""
1558
1559 # -----#
1560 # Celda 9: construcción de un payload JSON explicable para un expediente
1561 # -----#

```

```

1563 import json
1564
1565 def build_prediction_payload(
1566     predictor: PortablePredictor,
1567     df_raw: pd.DataFrame,
1568     row_index: int = 0,
1569     knn_topk: int = 10,
1570     knn_space: str = "all",
1571 ) -> Dict[str, Any]:
1572     """
1573         Construye un diccionario (JSON-serializable) con toda la información
1574         relevante para un único expediente.
1575
1576         Pensado para que un GPT (con Python) pueda:
1577             - ejecutar esta función,
1578             - obtener el dict,
1579             - y a partir de ahí generar una explicación en lenguaje natural.
1580
1581     Parámetros
1582     -----
1583     predictor : PortablePredictor
1584         Instancia ya cargada con load_predictor_from_zip().
1585     df_raw : DataFrame
1586         DataFrame con al menos una fila, en el formato esperado por predictor.
1587         (idealmente, df_norm tras export_python_transformacion).
1588     row_index : int
1589         Índice posicional (iloc) de la fila a explicar.
1590     knn_topk : int
1591         Número de vecinos más cercanos a devolver.
1592     knn_space : {"train", "all"}
1593         Conjunto de referencia para el kNN.
1594
1595     Devuelve
1596     -----
1597     dict
1598         Objeto listo para pasarse a json.dumps().
1599     """
1600     if not isinstance(df_raw, pd.DataFrame):
1601         raise TypeError("build_prediction_payload espera un pandas.DataFrame en"
1602                         "df_raw.")
1603     if df_raw.shape[0] == 0:
1604         raise ValueError("df_raw no contiene filas.")
1605
1606     if row_index < 0 or row_index >= df_raw.shape[0]:
1607         raise IndexError(f"row_index={row_index} fuera de rango para df_raw con {"
1608                         f"df_raw.shape[0]} filas.")
1609
1610     # Extraemos una sola fila como DataFrame (preserva nombres de columnas)
1611     df_one = df_raw.iloc[[row_index]].copy()
1612     idx_label = df_one.index[0]
1613
1614     # 1) Predicciones de todos los modelos
1615     res = predictor.predict(df_one)
1616
1617     ensemble_df = res["ensemble"]
1618     rf_base_out = res["rf_base"]
1619     rf_124_out = res["rf_124"]
1620     xgb_out = res["xgb_baja"]
1621
1622     # Fila única del ensemble
1623     ens_row = ensemble_df.loc[idx_label]
1624
1625     # RF base y RF_124: proba de la fila
1626     rf_base_proba_row = rf_base_out["proba"].loc[idx_label]
1627     rf_124_proba_row = rf_124_out["proba"].loc[idx_label]
1628
1629     # 2) Vecinos kNN para esa misma fila
1630     knn_res = predictor.knn_neighbors(df_one, topk=knn_topk, space=knn_space)
1631     knn_indices = knn_res["indices"][0].tolist()
1632     knn_distances = [float(d) for d in knn_res["distances"][0]]
1633     knn_similarity = [float(s) for s in knn_res["similarity_percent"][0]]

```

```

1633 p95_nn1_train = float(knn_res["p95_nn1_train"])
1634
1635 # 3) Entrada "cruda": valores de la fila en df_raw (ya normalizado)
1636 input_row = df_one.iloc[0].to_dict()
1637 # Convertimos posibles tipos numpy a tipos nativos
1638 input_row_clean = {}
1639 for k, v in input_row.items():
1640     if isinstance(v, (np.generic, np.bool_)):
1641         input_row_clean[k] = v.item()
1642     else:
1643         input_row_clean[k] = v
1644
1645 # 4) Construimos el payload
1646 payload = {
1647     "input": {
1648         "index_label": str(idx_label),
1649         "row_index": int(row_index),
1650         "features": input_row_clean,
1651     },
1652     "models": {
1653         "xgb_baja": {
1654             "p_bajo": float(xgb_out["p_bajo"].loc[idx_label]),
1655             "baja_comp": int(xgb_out["baja_comp"].loc[idx_label]),
1656             "threshold_bajo": float(xgb_out["threshold_bajo"]),
1657             "labels": list(predictor.xgb_baja.labels),
1658         },
1659         "rf_base": {
1660             "pred": str(rf_base_out["pred"].loc[idx_label]),
1661             "proba": {
1662                 str(cls): float(rf_base_proba_row[cls])
1663                 for cls in predictor.rf_base.classes_
1664             },
1665             "classes": [str(c) for c in predictor.rf_base.classes_],
1666         },
1667         "rf_124": {
1668             "pred": str(rf_124_out["pred"].loc[idx_label]),
1669             "proba": {
1670                 str(cls): float(rf_124_proba_row[cls])
1671                 for cls in predictor.rf_124.classes_
1672             },
1673             "classes": [str(c) for c in predictor.rf_124.classes_],
1674         },
1675     },
1676     "ensemble": {
1677         "pred_final": str(ens_row["pred"]),
1678         "confianza_pred": float(ens_row["confianza_pred"]),
1679         "fuente_confianza": str(ens_row["fuente_confianza"]),
1680         "flip_ensemble": bool(ens_row["flip_ensemble"]),
1681         "baja_comp": int(ens_row["baja_comp"]),
1682         "p_bajo": float(ens_row["p_bajo"]),
1683     },
1684     "knn": {
1685         "space": str(knn_space),
1686         "topk": int(knn_topk),
1687         "indices": [int(i) for i in knn_indices],
1688         "distances": knn_distances,
1689         "similarity_percent": knn_similarity,
1690         "p95_nn1_train": p95_nn1_train,
1691     },
1692     "meta": {
1693         "zip_path": str(predictor.zip_path),
1694         "notes": (
1695             "Probabilidades RF portables pueden diferir hasta ~0.02 respecto "
1696             "a sklearn; threshold_bajo del XGB se fija explícitamente en 0.38."
1697         ),
1698     },
1699 }
1700
1701 return payload
1702
1703
1704 # Ejemplo de uso rápido (puedes adaptarlo a tu df_norm real):

```

```

1705 #
1706 # zip_path =
1707 # predictor = load_predictor_from_zip(zip_path)
1708 #
1709 # import pandas as pd
1710 # df_test =
1711 pd.read_excel("/content/drive/MyDrive/_Pipeline_produccion/prediccion/modelos_portable.zip")
1712 # df_norm = df_test # sustituye aquí por df_norm si ya lo tienes transformado
1713 #
1714 # payload = build_prediction_payload(predictor, df_norm, row_index=0, knn_topk=10,
1715 knn_space="all")
1716 # print(json.dumps(payload, ensure_ascii=False, indent=2))
1717
1718 """Celda 11 - Documento RAG completo por expediente"""
1719
1720 # -----
1721 # Celda 11: documento RAG completo por expediente
1722 #           (JSON-safe + vecinos kNN + summary + risk_label + meta.version)
1723 # -----
1724
1725 import json
1726 from pathlib import Path
1727 import math
1728 import datetime
1729 from typing import Any, Dict
1730
1731 def _to_jsonable_scalar(v: Any) -> Any:
1732     """
1733         Convierte un valor escalar a un tipo JSON-serializable.
1734
1735         - np.generic / np.bool_      -> .item()
1736         - pandas.Timestamp          -> ISO8601 (str)
1737         - datetime.date/datetime   -> ISO8601 (str)
1738         - pandas.Timedelta          -> str (segundos)
1739         - float NaN                -> None (se convierte en null en JSON)
1740         - Resto                    -> se devuelve tal cual
1741     """
1742
1743     # Números y booleanos numpy (np.generic, np.bool_)
1744     if isinstance(v, (np.generic, np.bool_)):
1745         v = v.item()
1746
1747     # Fechas/horas de pandas
1748     if isinstance(v, pd.Timestamp):
1749         return v.isoformat()
1750
1751     # Fechas/horas estándar
1752     if isinstance(v, (datetime.datetime, datetime.date)):
1753         return v.isoformat()
1754
1755     # Duraciones
1756     if isinstance(v, pd.Timedelta):
1757         return str(v.total_seconds())
1758
1759     # NaN como null (JSON estándar)
1760     if isinstance(v, float) and math.isnan(v):
1761         return None
1762
1763     return v
1764
1765 def build_rag_case_document(
1766     predictor: PortablePredictor,
1767     df_raw: pd.DataFrame,
1768     row_index: int = 0,
1769     knn_topk: int = 10,
1770     knn_space: str = "all",
1771 ) -> Dict[str, Any]:
1772     """
1773         Construye un documento RAG v4 para un expediente concreto,
1774         listo para ser indexado o usado directamente por un GPT.

```

```

1773
1774 Estructura principal:
1775   - doc_type : "licitacion_riesgo_prediccion"
1776   - rag_metodologia : info sobre la metodología RAG v4 (Source, etc.)
1777   - source : payload factual (predicciones, kNN, entrada)
1778   - retriever : campos clave para filtrado/búsqueda en el índice
1779   - evaluation : placeholders para métricas/feedback futuro
1780
1781 Además:
1782   - 'source.knn.neighbors' contiene ya los datos completos de los casos
1783     vecinos (filas de df_raw), incluyendo un bloque 'summary' con las
1784     6 variables principales:
1785       * Ahorro_final
1786       * Presupuesto_licitacion_lote_c
1787       * N_ofertantes
1788       * Plazo_m
1789       * Baja_p
1790       * C_precio_p
1791   - 'source.models.xgb_baja' incorpora:
1792     * 'proba': probabilidades por clase (BAJO/MEDIO/ALTO)
1793     * 'risk_label': etiqueta final de riesgo de baja competencia
1794   - 'source.meta' incorpora:
1795     * 'model_set_version'
1796     * 'generated_at' (UTC, ISO8601)
1797 """
1798 # 1) Payload factual base (una fila + kNN)
1799 payload = build_prediction_payload(
1800     predictor=predictor,
1801     df_raw=df_raw,
1802     row_index=row_index,
1803     knn_topk=knn_topk,
1804     knn_space=knn_space,
1805 )
1806
1807 # --- Meta: versión de modelos y timestamp de generación ---
1808 meta_payload = payload.get("meta", {}) or {}
1809 meta_payload["model_set_version"] = "v1.0_portable_rag"
1810 meta_payload["generated_at"] = datetime.datetime.utcnow().isoformat() + "Z"
1811 payload["meta"] = meta_payload
1812
1813 # --- Riesgo XGB: proba completa + risk_label ---
1814 xgb_block = payload["models"]["xgb_baja"]
1815 # Recalculamos las probabilidades completas para las 3 clases del XGB
1816 df_one = df_raw.iloc[[row_index]].copy()
1817 proba_all = predictor.xgb_baja.predict_proba(df_one) # shape (1, num_class)
1818 proba_vec = proba_all[0]
1819 labels_list = list(predictor.xgb_baja.labels_)
1820
1821 # Lista python de probabilidades para trabajar cómodamente
1822 proba_list = [float(p) for p in proba_vec]
1823 # Guardamos el vector completo en el bloque XGB
1824 xgb_block["proba"] = {
1825     str(lbl): float(proba_list[i]) for i, lbl in enumerate(labels_list)
1826 }
1827
1828 # Definimos risk_label:
1829 #   - Si baja_comp = 1 y existe clase "BAJO" -> "BAJO"
1830 #   - Si baja_comp = 0 -> clase con mayor probabilidad entre el resto
1831 #     (MEDIO/ALTO, etc.); si no hay "BAJO", usamos argmax global.
1832 baja_flag = int(xgb_block["baja_comp"])
1833 if "BAJO" in labels_list:
1834     idx_bajo = labels_list.index("BAJO")
1835 else:
1836     idx_bajo = None
1837
1838 if baja_flag == 1 and idx_bajo is not None:
1839     risk_label = "BAJO"
1840 else:
1841     if idx_bajo is not None and len(labels_list) > 1:
1842         # Candidatas: todas menos BAJO
1843         cand = [
1844             (lbl, proba_list[i])

```

```

1845         for i, lbl in enumerate(labels_list)
1846             if i != idx_bajo
1847         ]
1848     if cand:
1849         risk_label = max(cand, key=lambda t: t[1])[0]
1850     else:
1851         # Fallback: argmax global
1852         idx_max = max(range(len(proba_list)), key=lambda i: proba_list[i])
1853         risk_label = labels_list[idx_max]
1854     else:
1855         # Fallback: argmax global
1856         idx_max = max(range(len(proba_list)), key=lambda i: proba_list[i])
1857         risk_label = labels_list[idx_max]
1858
1859     xgb_block["risk_label"] = str(risk_label)
1860     payload["models"]["xgb_baja"] = xgb_block
1861
1862     # --- Limpieza JSON-safe de las features de entrada ---
1863     features_raw = payload["input"]["features"]
1864     features_clean = {k: _to_jsonable_scalar(v) for k, v in features_raw.items()}
1865     payload["input"]["features"] = features_clean
1866
1867     # Identificadores básicos del caso
1868     index_label = payload["input"]["index_label"]
1869     row_idx = payload["input"]["row_index"]
1870     ensemble = payload["ensemble"]
1871     rf_base = payload["models"]["rf_base"]
1872     rf_124 = payload["models"]["rf_124"]
1873
1874     # 2) Enriquecer la sección kNN con datos completos de vecinos + summary
1875     knn_section = payload.get("knn", {}) or {}
1876     indices = knn_section.get("indices", []) or []
1877     neighbors_list = []
1878     n_rows = df_raw.shape[0]
1879
1880     for rank, idx in enumerate(indices):
1881         # Índice en Z_all (y, asumimos, en df_raw)
1882         try:
1883             idx_int = int(idx)
1884         except Exception:
1885             continue
1886
1887         if idx_int < 0 or idx_int >= n_rows:
1888             # Por robustez, si el índice está fuera de rango
1889             neighbors_list.append(
1890                 {
1891                     "rank": int(rank + 1),
1892                     "knn_index": idx_int,
1893                     "case_id": None,
1894                     "summary": None,
1895                     "features": None,
1896                 }
1897             )
1898         continue
1899
1900
1901     # Fila completa del vecino en df_raw
1902     row = df_raw.iloc[idx_int]
1903     row_dict = row.to_dict()
1904     row_clean = {k: _to_jsonable_scalar(v) for k, v in row_dict.items()}
1905
1906     # Bloque summary con las 6 variables principales
1907     summary = {
1908         "Ahorro_final": row_clean.get("Ahorro_final"),
1909         "Presupuesto_licitacion_lote_c": row_clean.get(
1910             "Presupuesto_licitacion_lote_c"),
1911         "N_ofertantes": row_clean.get("N_ofertantes"),
1912         "Plazo_m": row_clean.get("Plazo_m"),
1913         "Baja_p": row_clean.get("Baja_p"),
1914         "C_precio_p": row_clean.get("C_precio_p"),
1915     }
1916
1917     neighbor_entry = {

```

```

1916     "rank": int(rank + 1),
1917     "knn_index": idx_int,
1918     # Intentamos usar el Identificador PLACSP como ID de caso; si no está,
1919     # usamos el índice
1920     "case_id": _to_jsonable_scalar(row_dict.get("Identificador", idx_int)),
1921     "summary": summary,
1922     "features": row_clean,
1923   }
1924   neighbors_list.append(neighbor_entry)
1925
1926   # Insertamos los vecinos en la sección kNN
1927   knn_section["neighbors"] = neighbors_list
1928   payload["knn"] = knn_section
1929
1930   # 3) Bloque retriever: campos útiles para indexar y filtrar en el RAG
1931   retriever = {
1932     # Identificador primario (índice del DataFrame)
1933     "case_id": index_label,
1934     "row_index": row_idx,
1935     # Identificadores "de negocio" (PLACSP) si existen
1936     "business_ids": {
1937       "Identificador": features_clean.get("Identificador"),
1938       "Identificador_lote": features_clean.get("Identificador_lote"),
1939       "Identificador_lici": features_clean.get("Identificador_lici_c"),
1940     },
1941     # Etiquetas de salida clave
1942     "labels": {
1943       "pred_ensemble": ensemble["pred_final"],
1944       "baja_comp_flag": ensemble["baja_comp"],
1945       "p_bajo": ensemble["p_bajo"],
1946       "rf_base_pred": rf_base["pred"],
1947       "rf_124_pred": rf_124["pred"],
1948       "xgb_risk_label": xgb_block["risk_label"],
1949     },
1950     # Algunas features clave para filtrado rápido
1951     "key_features": {
1952       k: v
1953       for k, v in features_clean.items()
1954       if k in {
1955         "Presupuesto_licitacion_lote_c",
1956         "N_ofertantes",
1957         "Plazo_m",
1958         "Baja_p",
1959         "C_precio_p",
1960         "C_juicio_valor_p_c",
1961         "Mes_lici",
1962         "Tipo_de_Administracion_c",
1963         "Tipo_de_procedimiento_c",
1964         "Provincia2",
1965         "CPV_c",
1966         "CPV_lote_c",
1967       }
1968     },
1969     # Lista completa de campos disponibles para filtrado/consulta avanzada
1970     "available_feature_keys": list(features_clean.keys()),
1971   }
1972
1973   # 4) Metadatos de metodología RAG v4 (informativos)
1974   rag_metodologia = {
1975     "version": "v4",
1976     "modules": [
1977       {
1978         "name": "Source",
1979         "role": (
1980           "Preparar y procesar datos estructurados de predicción, "
1981           "incluyendo entradas, salidas de modelos y casos similares."
1982         ),
1983       },
1984       {
1985         "name": "UI",
1986         "role": (
1987           "Interfaz de usuario que mostrará al licitador la explicación"
1988         )
1989       }
1990     ]
1991   }

```

```

1987         "basada en este documento (fuera de este script)."
1988     ),
1989 },
1990 {
1991     "name": "Orchestration",
1992     "role": (
1993         "Coordinar la carga de este JSON, la invocación al GPT y el "
1994         "almacenamiento de feedback (gestionado fuera de este script)."
1995     ),
1996 },
1997 {
1998     "name": "Retriever",
1999     "role": (
2000         "Recuperar documentos relevantes usando las claves de "
2001         "'retriever', "
2002         "por ejemplo por pred_ensemble, baja_comp_flag o CPV."
2003     ),
2004 },
2005 {
2006     "name": "Generator",
2007     "role": (
2008         "Generar la explicación en lenguaje natural a partir de este "
2009         "documento y otras fuentes (tesis, normativa, etc.)."
2010     ),
2011 },
2012 {
2013     "name": "Reranker",
2014     "role": (
2015         "Reordenar documentos recuperados si hay múltiples casos "
2016         "candidatos para una consulta dada."
2017     ),
2018 },
2019 {
2020     "name": "Evaluation",
2021     "role": (
2022         "Medir y registrar la calidad de las respuestas generadas "
2023         "usando este documento como base factual."
2024     ),
2025 ],
2026 }
2027
2028 # 5) Placeholder para evaluación (se rellenará en otro módulo del sistema)
2029 evaluation = {
2030     "human_feedback": None,
2031     "notes": None,
2032     "metrics": {
2033         "retrieval": None,
2034         "generation": None,
2035     },
2036 }
2037
2038 # 6) Documento RAG completo
2039 rag_doc = {
2040     "doc_type": "licitacion_riesgo_prediccion",
2041     "rag_metodologia": rag_metodologia,
2042     "source": payload,           # payload ya con features, vecinos y XGB
2043     "retriever": retriever,
2044     "evaluation": evaluation,
2045 }
2046
2047 return rag_doc
2048
2049
2050 def generate_and_save_rag_case_document(
2051     predictor: PortablePredictor,
2052     df_raw: pd.DataFrame,
2053     row_index: int,
2054     knn_topk: int = 10,
2055     knn_space: str = "all",
2056     output_path: str = "rag_case_doc.json",

```

```
2057 ) -> Dict[str, Any]:  
2058     """  
2059     Genera el documento RAG completo para un expediente y lo guarda en JSON.  
2060  
2061     Este será el archivo 'legible y utilizable' por el GPT:  
2062     - Se puede cargar con json.load(...)  
2063     - El GPT lo usará como documento fuente en el RAG, sin necesidad de  
2064       acceder a df_norm ni a otros artefactos de kNN.  
2065     """  
2066     rag_doc = build_rag_case_document(  
2067         predictor=predictor,  
2068         df_raw=df_raw,  
2069         row_index=row_index,  
2070         knn_topk=knn_topk,  
2071         knn_space=knn_space,  
2072     )  
2073  
2074     output_path = Path(output_path)  
2075     with output_path.open("w", encoding="utf-8") as f:  
2076         json.dump(rag_doc, f, ensure_ascii=False, indent=2)  
2077  
2078     print(f"Documento RAG guardado en {output_path.resolve() }")  
2079     return rag_doc
```