

```

1 # Guía técnica para el GPT - Uso de `gpt_portable_predictor.py` y documentos RAG de
2 predicción
3
4 Este documento describe cómo debe trabajar el GPT con:
5
6 - El **motor de predicción portable** implementado en `gpt_portable_predictor.py`.
7 - Los artefactos de modelos contenidos en `modelos_portable.zip`.
8 - Los documentos JSON de salida tipo `rag_case_doc_ejemplo.json`.
9
10 Su objetivo es que el GPT pueda:
11
12 1. Recoger datos de una licitación (a partir de diálogo con el usuario).
13 2. Normalizarlos a un formato tabular compatible con el motor.
14 3. Ejecutar el motor en Python para obtener un documento RAG por expediente.
15 4. Explicar los resultados al usuario, apoyándose en dicho JSON y en la documentación
16 adicional (tesis, metodología, normativa).
17
18 ---  

19
20 ## 1. Flujo general del sistema
21
22 El sistema completo funciona en dos capas:
23
24 1. **Motor de predicción portable (Python)**
25 - Ficheros principales:
26   - `gpt_portable_predictor.py`
27   - `modelos_portable.zip`
28 - Se ejecuta en el entorno Python del GPT.
29 - Produce, para cada expediente, un único objeto JSON estructurado (`rag_case_doc`)
30   con:
31     - Entrada normalizada.
32     - Salidas de modelos (RF, XGB).
33     - Predicción final del ensemble.
34     - Vecinos kNN con información detallada.
35     - Metadatos y campos pensados para RAG.
36
37 2. **GPT-RAG orquestador**
38 - Mantiene una conversación con el usuario.
39 - Traduce las respuestas del usuario a las variables que requiere el motor.
40 - Llama a Python para generar el `rag_case_doc`.
41 - Usa ese JSON como **fuente factual principal** para explicar el resultado,
42 combinándolo con:
43   - Metodología RAG.
44   - Capítulos de tesis.
45   - Normativa relacionada, si es relevante.
46
47 El GPT **no debe re-entrenar modelos ni modificar artefactos**. Únicamente debe:
48
49 - Preparar inputs,
50 - Invocar el motor,
51 - Interpretar y explicar las salidas.
52
53 ---  

54
55 ## 2. Motor de predicción portable (`gpt_portable_predictor.py`)
56
57 ### 2.1. Carga de modelos
58
59 El punto de entrada es:
60
61
62
63   - `modelos_portable.zip` contiene:
64     - `preprocess_pre_portable.json`
65     - `rf_base_portable.json`
66     - `rf_124_portable.json`
67     - `xgb_baja_portable.json`
68     - `xgb_baja_model.json`
```

```

69     - `knn_space_portable.json`
70     - `knn_train.npz`
71     - `knn_all.npz`
72     - `ensemble_portable.json`
73     - `centroides_pre.json`
74 - El motor **no utiliza sklearn ni xgboost**. Reimplementa:
75     - Árboles de RandomForest desde JSON.
76     - Booster de XGBoost desde JSON.
77     - Preprocesado (imputación + escalado + one-hot encoding).
78     - Búsqueda de vecinos kNN en espacio PRE_NUM.
79
80 ### 2.2. Clases y funciones principales
81
82 Las piezas internas se organizan en clases (no es necesario usarlas directamente desde el GPT, salvo a través de la fachada):
83
84     - `PortablePreprocessor`
85     - `RandomForestPortable`
86     - `XGBBinaryPortable`
87     - `KNNPortable`
88     - `EnsemblePredictor`
89     - `PortablePredictor` (fachada principal)
90
91 La interfaz pública relevante es:
92
93 ```python
94 from gpt_portable_predictor import (
95     load_predictor_from_zip,
96     build_prediction_payload,
97     build_rag_case_document,
98     generate_and_save_rag_case_document,
99 )
100 ```

101 #### 2.2.1. `predictor.predict(df_raw)`
102
103 ```python
104 resultados = predictor.predict(df_raw)
105 ```

106 Devuelve un diccionario con al menos:
107
108
109     - `"rf_base"`:
110         - `pred`: Series con clase predicha (cluster) para cada fila.
111         - `proba`: DataFrame con probabilidades por clase.
112     - `"rf_124"`:
113         - Igual estructura que `rf_base` pero para el modelo especializado en clases 1,2,4.
114     - `"xgb_baja"`:
115         - `p_bajo`: probabilidad continua de "baja competencia" según XGB.
116         - `baja_comp`: flag binario (1 si `p_bajo >= 0.38`, 0 en caso contrario).
117         - `threshold_bajo`: umbral **fijado explícitamente en 0.38**.
118         - `labels`: lista de etiquetas de salida del XGB (ej. `["BAJO","MEDIO","ALTO"]`).
119
120 #### 2.2.2. `predictor.knn_neighbors(df_raw, topk, space)`
121
122
123 ```python
124 vecinos = predictor.knn_neighbors(df_raw, topk=10, space="all")
125 ```

126 Devuelve:
127
128
129     - `indices`: array de índices de las filas vecinas (en `z_train` o `z_all`).
130     - `distances`: array de distancias euclídeas en el espacio PRE_NUM.
131     - `similarity_percent`: similitud relativa en [%], derivada de la distancia y un valor de referencia `p95_nn1_train`.
132     - `p95_nn1_train`: percentil 95 de la distancia al primer vecino en TRAIN.
133
134 > **Nota**: el script de generación de RAG ya "resuelve" estos índices y rellena los datos completos de los vecinos; el GPT suele trabajar directamente con el JSON RAG y no con esta función.
135

```

```

136 ##### 2.2.3. Funciones de alto nivel para RAG
137
138 - `build_prediction_payload(predictor, df_raw, row_index, knn_topk, knn_space)`
  Construye un payload estructurado (dict JSON-safe) con:
  - Entrada (`input`),
  - Salidas de modelos (`models`),
  - Ensemble (`ensemble`),
  - kNN (`knn`),
  - Metadatos (`meta`).
145
146 - `build_rag_case_document(predictor, df_raw, row_index, knn_topk, knn_space)`
  Construye el documento RAG completo (dict).
148
149 - `generate_and_save_rag_case_document(...)`
  Igual que la anterior, pero además lo guarda en un `json` en disco.
151
152 ---
153
154 ## 3. Contrato de entrada para nuevas predicciones
155
156 ### 3.1. Estructura general del input
157
158 Para un nuevo expediente, el GPT debe construir un `DataFrame` de una sola fila con
  **nombres de columna compatibles** con el modelo (los mismos que en
  `df_test_train_v8.xlsx` / `df_norm`).
159
160 Es suficiente que el GPT construya un `dict` con las Variables de **PRE-licitación**
  (numéricas y categóricas):
161
162 ````python
163 user_input = {
164     "Presupuesto_licitacion_lote_c": ....,
165     "Plazo_m": ....,
166     "C_precio_p": ....,
167     "C_juicio_valor_p_c": ....,
168     "Mes_lici": ....,
169     "Tipo_de_Administracion_c": ....,
170     "Tipo_de_procedimiento_c": ....,
171     "N_lotes_c": ....,
172     "N_clasi_empresa_C": ....,
173     "Tipo_de_contrato_c": ....,
174     "Codigo_Postal_c": ....,
175     "N_CPV_c": ....,
176 }
177 df_norm = pd.DataFrame([user_input])
178 ````

179 Ese `df_norm` se pasa al motor (`predictor` y funciones RAG).
180
181
182 ### 3.2. Variables clave (las más importantes)
183
184 El GPT debería priorizar la recogida y normalización correcta de, al menos, estas
  variables:
185
186 | Variable                                     | Tipo          | Unidad / dominio           | Comentario
187 | Variable                                     | Tipo          | Unidad / dominio           | Comentario
188 | `Presupuesto_licitacion_lote_c`             | float         | euros                      | Presupuesto
  de licitación del lote. Siempre en euros, valor positivo.
189 | `Plazo_m`                                    | int           | meses                     | Plazo
  contractual en meses.
190 | `C_precio_p`                                | float         | % en 0-100                | Ponderación
  del precio en la valoración.
191 | `C_juicio_valor_p_c`                         | float         | % en 0-100                | Ponderación de
  criterios de juicio de valor.

```

```

192 | `Mes_lici`           | int    | 1-12          | Mes de
licitación (enero=1,..., diciembre=12).
193 | `Tipo_de_Administracion_c` | str    | categoría      | Tipo de
administración (ej.: `"Autoridad local"`, `"Administración General"`, etc.).
194 | `Tipo_de_procedimiento_c` | str    | categoría      | Procedimiento
de contratación (ej.: `"Abierto"`, `"Abierto simplificado"`, etc.).
195 | `Provincia2`          | int    | código numérico | Código de
provincia obtenido del código postal.
196 | `N_CPV`                | int    | unidades enteras | Número de
clasificadores CPV de la licitación.
197 | `N_lotes`               | int    | unidades enteras | Número de lotes en
los que se divide la licitación.
198 | `N_clasi_empresa`     | int    | unidades enteras | Número de
requerimientos empresariales exigidos en la licitación.
199 | `Tipo_de_contrato_c`   | str    | categoría      | Tipo de contrato
(ej.: `"Obras"`, `"Servicios"`, etc.).
200
201 Además, existen otras variables (por ejemplo, fechas de publicación, identificadores
PLACSP, etc.) que pueden completarse si el usuario las proporciona, pero no son
imprescindibles para la predicción.

```

3.3. Reglas de normalización desde lenguaje natural

El GPT debe aplicar reglas claras para transformar entradas textuales del usuario en valores numéricos/categóricos:

- Meses:
 - `"Octubre"`, `"octubre (10)"`, `"mes 10"` → `Mes_lici = 10`.
- Porcentajes:
 - `"75%"`, `"75 %"` → `75.0`.
 - `"0.75"` ****interpretar preferentemente como 75 %****, salvo que el contexto indique explícitamente 0-1.
- Plazos:
 - `"14 meses"`, `"14 m"`, `"duración de 14"` → `Plazo_m = 14`.
- Presupuesto:
 - `"893.145 €"`, `"893145"`, `"0,893 M€"` → `Presupuesto_licitacion_lote_c ≈ 893145.0`
 - El GPT debe intentar entender unidades (miles, millones) y expresarlo siempre en euros.

Si la interpretación es ambigua, el GPT debe:

- Pedir aclaración, o
- Explicitar la suposición que va a realizar y advertir de la limitación.

3.4. Gestión de valores faltantes

Si el usuario no puede proporcionar alguna variable:

- El motor seguirá funcionando gracias a la imputación interna (medias, modas).
- Sin embargo, desde el punto de vista del DSS, la predicción será menos fiable.

El GPT debe:

- Indicar al usuario qué campos se han imputado o dejado por defecto.
- Señalar explícitamente que esto reduce la confianza en la interpretación, aunque el motor devuelva un resultado numérico.

4. Estructura del documento RAG (`rag_case_doc`)

4.1. Esquema de alto nivel

El documento JSON generado por `build_rag_case_document` tiene esta estructura:

```

242
243 ```json
244 {
245   "doc_type": "licitacion_riesgo_prediccion",
246   "rag_metodologia": { ... },
247   "source": {
248     "input": { ... },
249     "models": { ... },
250     "ensemble": { ... },
251     "knn": { ... },
252     "meta": { ... }
253   },
254   "retriever": { ... },
255   "evaluation": { ... }
256 }
257 ```
258
259 ### 4.2. `source.input`
260
261 Ejemplo:
262
263 ```json
264 "input": {
265   "index_label": "0",
266   "row_index": 0,
267   "features": {
268     "Identificador": 4353889,
269     "Presupuesto_licitacion_lote_c": 520161.64,
270     "Plazo_m": 4,
271     "C_precio_p": 100.0,
272     "C_juicio_valor_p_c": 0.0,
273     "Mes_lici": 12,
274     "Tipo_de_Administracion_c": "Autoridad local",
275     "Tipo_de_procedimiento_c": "Abierto simplificado",
276     "N_CPV": "3",
277     ...
278   }
279 }
280 ```
281
282 - `features` contiene **todas las columnas** del `df_norm` para ese caso.
283 - Fechas se representan como cadenas ISO (`"YYYY-MM-DDTHH:MM:SS"`).
284 - Valores faltantes se representan como `null` (JSON estándar).
285
286 ### 4.3. `source.models`
287
288 El bloque `models` contiene, al menos:
289
290 #### 4.3.1. `xgb_baja`
291
292 ```json
293 "xgb_baja": {
294   "p_bajo": 0.4697,
295   "baja_comp": 1,
296   "threshold_bajo": 0.38,
297   "labels": ["BAJO", "MEDIO", "ALTO"],
298   "proba": {
299     "BAJO": 0.4697,
300     "MEDIO": 0.4175,
301     "ALTO": 0.1126
302   },
303   "risk_label": "BAJO"
304 }
305 ```
306
307 Interpretación:
308
309 - `p_bajo`: probabilidad continua de "baja competencia" (clase BAJO).
310 - `baja_comp`: flag binario, calculado como `p_bajo >= 0.38`.
311 - `proba`: distribución completa entre BAJO / MEDIO / ALTO.
312 - `risk_label`:
313   - `BAJO` cuando `baja_comp = 1` y existe etiqueta BAJO.

```

```

314     - Si `baja_comp = 0`, es la etiqueta más probable entre el resto de clases
315     (típicamente MEDIO o ALTO).
316
317 #### 4.3.2. `rf_base` y `rf_124`
318
319 Ambos tienen estructura similar, por ejemplo:
320
321 ```json
322 "rf_base": {
323     "pred": "4",
324     "proba": {
325         "0": 0.037,
326         "1": 0.072,
327         "2": 0.068,
328         "3": 0.172,
329         "4": 0.496,
330         "5": 0.154
331     },
332     "classes": ["0", "1", "2", "3", "4", "5"]
333 }
334 ````
335     - `rf_base`: modelo general de clasificación en 6 clusters.
336     - `rf_124`: modelo especializado que sólo discrimina entre clases 1, 2 y 4.
337
338 Las probabilidades del RF portable pueden diferir hasta ~0.02 respecto a sklearn.
339 Este margen es aceptado y está documentado en `meta.notes`.
340
341 #### 4.4. `source.ensemble`
342
343 Ejemplo:
344
345 ```json
346 "ensemble": {
347     "pred_final": "1",
348     "confianza_pred": 0.7151,
349     "fuente_confianza": "rf_124_basebc",
350     "flip_ensemble": true,
351     "baja_comp": 1,
352     "p_bajo": 0.4697
353 }
354 ````
355
356 Interpretación:
357
358     - `pred_final`: clase de cluster final (texto `"0"`-`"5"`).
359     - `confianza_pred`: probabilidad de la clase ganadora en el modelo que se ha usado
360 como fuente de decisión.
361     - `fuente_confianza`:
362         - `rf_base` o `rf_124_basebc`, indica qué modelo fue elegido.
363     - `flip_ensemble`:
364         - `true` indica que se ha cambiado la predicción de RF base por la de RF_124.
365     - `baja_comp`, `p_bajo`:
366         - Copia de los valores de XGB para integrarlos en la explicación global.
367
368 Regla de combinación implementada:
369
370 1. Se ejecutan RF base y RF_124.
371 2. Si RF base predice una clase **dentro de {1,2,4}**, se compara la probabilidad
372 top-1 de RF base y RF_124.
373 3. Se elige el modelo con mayor probabilidad top-1:
374     - Si se elige RF_124 → `flip_ensemble = true`.
375     - Si se mantiene RF base → `flip_ensemble = false`.
376
377 #### 4.5. `source.knn`
378
379 Ejemplo de estructura:
380
381 ```json
382 "knn": {
383     "space": "all",
384     "topk": 10,

```

```

383     "indices": [0, 999, 49, ...],
384     "distances": [0.0, 0.0168, ...],
385     "similarity_percent": [100.0, 98.9, ...],
386     "p95_nn1_train": 1.5866,
387     "neighbors": [
388         {
389             "rank": 1,
390             "knn_index": 0,
391             "case_id": 4353889,
392             "summary": {
393                 "Ahorro_final": 9.72,
394                 "Presupuesto_licitacion_lote_c": 520161.64,
395                 "N_ofertantes": 8,
396                 "Plazo_m": 4,
397                 "Baja_p": 20.4,
398                 "C_precio_p": 100.0
399             },
400             "features": { ... fila completa del vecino ... }
401         },
402         ...
403     ]
404 }
405 ```
406
407 - `neighbors` ya incluye:
408   - `summary` con las 6 variables más relevantes (ahorro y 5 económicas).
409   - `features` con la fila completa del `df_norm` del vecino.
410 - El GPT **no necesita acceder a ningún otro dataset** para explicar los casos parecidos; toda la información está en este bloque.
411
412 ### 4.6. `source.meta`
413
414 Contiene, al menos:
415
416 ```json
417 "meta": {
418     "zip_path": "ruta/relativa/a/modelos_portable.zip",
419     "notes": "Probabilidades RF portables pueden diferir hasta ~0.02 respecto a sklearn; threshold_bajo del XGB se fija explícitamente en 0.38.",
420     "model_set_version": "v1.0_portable_rag",
421     "generated_at": "2025-11-16T15:45:04.240573Z"
422 }
423 ```
424
425 - `model_set_version`: etiqueta de versión de artefactos.
426 - `generated_at`: timestamp de inferencia en UTC (ISO8601, con sufijo `Z`).
427
428 ### 4.7. `retriever` y `evaluation`
429
430 - `retriever` agrupa claves útiles para indexado y filtrado:
431
432 ```json
433 "retriever": {
434     "case_id": "0",
435     "row_index": 0,
436     "business_ids": {
437         "Identificador": 4353889,
438         "Identificador_lote": "4353889//NEGOBP-2019/1200",
439         "Identificador_lici": "NEGOBP-2019/1200"
440     },
441     "labels": {
442         "pred_ensemble": "1",
443         "baja_comp_flag": 1,
444         "p_bajo": 0.4697,
445         "rf_base_pred": "4",
446         "rf_124_pred": "1",
447         "xgb_risk_label": "BAJO"
448     },
449     "key_features": {
450         "Presupuesto_licitacion_lote_c": 520161.64,
451         "N_ofertantes": 8,
452         "Plazo_m": 4,

```

```

453     "Baja_p": 20.4,
454     "C_precio_p": 100.0,
455     "C_juicio_valor_p_c": 0.0,
456     "Mes_lici": 12,
457     "Tipo_de_Administracion_c": "Autoridad local",
458     "Tipo_de_procedimiento_c": "Abierto simplificado",
459     "Provincia2": 33,
460     "CPV_c": "45233252-0",
461     "CPV_lote_c": "45233252-0"
462   },
463   "available_feature_keys": [ ... lista completa de campos ... ]
464 }
```
466
467 - `evaluation` es un placeholder para métricas y feedback, que puede rellenarse en etapas posteriores.
468
469 ---
470
471 ## 5. Pautas para el GPT al usar este sistema
472
473 ### 5.1. Cuándo ejecutar Python
474
475 El GPT debe ejecutar código Python cuando:
476
477 - El usuario solicite **una nueva predicción** ("quiero estimar el riesgo de esta licitación...").
478 - El usuario quiera **recalcular** o **analizar un caso específico** con nuevos datos.
479
480 Secuencia típica:
481
482 1. Conversación para recoger datos.
483 2. Normalización a un `dict` con las variables clave.
484 3. En Python:
485 ``python
486 from gpt_portable_predictor import load_predictor_from_zip,
487 generate_and_save_rag_case_document
488 import pandas as pd
489
490 predictor = load_predictor_from_zip("modelos_portable.zip")
491 df_norm = pd.DataFrame([user_input_dict])
492 rag_doc = generate_and_save_rag_case_document(
493 predictor=predictor,
494 df_raw=df_norm,
495 row_index=0,
496 knn_topk=10,
497 knn_space="all",
498 output_path="rag_case_doc_tmp.json",
499)
```
500 4. De vuelta al usuario: explicación basada en `rag_doc`.
501
502 ### 5.2. Cómo explicar resultados
503
504 El GPT debe:
505
506 - Basarse **primariamente** en `rag_doc["source"]` para:
507   - Clase final (`ensemble.pred_final`) y confianza (`ensemble.confianza_pred`).
508   - Riesgo de baja competencia (`xgb_baja.p_bajo`, `xgb_baja.risk_label`, `baja_comp`).
509   - Comparación con vecinos (`knn.neighbors[i].summary` y `features`).
510 - Usar la tesis/metodología como **contexto**:
511   - para interpretar clusters,
512   - para justificar por qué determinadas variables son relevantes,
513   - para conectar con conceptos estadísticos y de gestión de proyectos.
514
515 Debe dejar claro al usuario:
516
517 - Qué partes de la explicación se basan en datos concretos del JSON.
518 - Qué partes son interpretación o contexto general (tesis, normativa).
519
520 ### 5.3. Qué hacer si la entrada es deficiente o incoherente

```

521
522 Si, tras conversar, el GPT detecta:
523
524 - Campos críticos sin valor (p. ej. presupuesto, n° de ofertas, baja).
525 - Valores claramente fuera de rango (p. ej. `Baja_p = -5`, `N_ofertantes = 0`).
526
527 Debe:
528
529 1. Intentar aclarar con el usuario.
530 2. Si no es posible, explicitar:
531 - Qué valores se han imputado o ajustado.
532 - Que la fiabilidad de la predicción se ve afectada.
533
534 El motor ****no lanza errores en estos casos****; simplemente imputa. El GPT es
responsable de advertir de esta situación.
535
536 **### 5.4. Errores que debe evitar el GPT**
537
538 - No debe:
539 - Re-entrenar modelos.
540 - Modificar los archivos del ZIP.
541 - Inventar columnas que el motor no usa.
542 - No debe asumir:
543 - Igualdad numérica exacta con sklearn/xgboost originales.
544 - Hay una tolerancia aceptada de ≈0,02 en probabilidades de RF.
545 - No debe reinterpretar:
546 - Los clusters como categorías que no estén respaldadas por la tesis/metodología.
547 - El riesgo de baja competencia sin respetar `risk_label` y `threshold_bajo`.
548
549 ---
550
551 Este documento debe usarse como ****guía interna del GPT**** para interactuar
correctamente con el motor portable y los documentos RAG, garantizando:
552
553 - Coherencia en la construcción de inputs.
554 - Interpretación fiel de los outputs.
555 - Transparencia con el usuario sobre supuestos, límites y nivel de confianza.
556
557
558 ---
559 **## 7. Modo light sin Python (usuarios sin herramienta de código)**
560
561 En algunos entornos o tipos de cuenta el GPT ****no puede ejecutar código Python****. En
ese escenario:
562
563 - No es posible usar `gpt_portable_predictor.py` ni cargar `modelos_portable.zip`.
564 - No se pueden recalcular distancias kNN ni probabilidades RF/XGB.
565 - El sistema solo puede trabajar con la ****información ya precomputada**** y cargada como
conocimiento estático.
566
567 Para estos casos se ha preparado un archivo adicional:
568
569 - `casos_referencia_light.md`
570
571 Este archivo contiene, para todos los casos TRAIN+TEST de la base utilizada en el DSS:
572
573 - Variables de ****PRE-licitación**** (numéricas y categóricas), incluyendo, entre otras:
574 - `Presupuesto_licitacion_lote_c`
575 - `C_precio_p`
576 - `Plazo_m`
577 - `N_lotes`
578 - `N_CPV`
579 - `N_clasi_empresa`
580 - `C_juicio_valor_p_c`
581 - `Intervalo_lici_d_c`
582 - `Tipo_de_contrato_c`
583 - `Tipo_de_Administracion_c`
584 - `Tipo_de_procedimiento_c`
585 - `Codigo_Postal_c`
586 - Las cinco variables económico-competitivas principales:
587 - `Presupuesto_licitacion_lote_c`
588 - `N_ofertantes`

```

589     - `Plazo_m`
590     - `Baja_p`
591     - `C_precio_p`
592   - El resultado histórico para cada caso:
593     - `Ahorro_final`
594     - `Cluster_6` (en el campo `cluster_final`)
595   - Coordenadas normalizadas `z_*` para las cinco variables principales, calculadas
      sobre toda la base histórica.

596
597 Cada caso se presenta como un bloque `json` con la siguiente información mínima:
598
599   - `case_id`: identificador del caso (normalmente `Identificador` de la BD).
600   - `dataset_index`: índice de fila en el dataset original.
601   - `cluster_final`: valor de `Cluster_6`.
602   - `Ahorro_final`.
603   - Las cinco variables principales y sus `z_*`.
604   - Un diccionario `pre_variables` con todas las variables de PRElicitación disponibles.

605
606 ### 7.1. Qué debe hacer el GPT en modo light
607
608 Cuando no haya Python disponible, el GPT debe seguir este enfoque:
609
610 1. **Recoger variables clave del usuario**
611 Al menos:
612   - `Presupuesto_licitacion_lote_c` (euros),
613   - `Plazo_m`,
614   - `C_precio_p`,
615 usando las mismas reglas de normalización que en el modo completo (porcentajes
0-100, meses 1-12, etc.).
616 Siempre que sea posible, también debería preguntar por otras variables de
PRElicitación (tipo de contrato, administración, procedimiento, código postal,
etc.).
```

617

618 **2. **Formular una consulta RAG sobre `casos_referencia_light.md`****

619 La consulta debe reflejar explícitamente los valores o rangos de:

- `Presupuesto_licitacion_lote_c`,
- `C_precio_p`,
- `Plazo_m`,
- y, si se conocen, el resto de PRE, ya que el RAG trabaja sobre texto.

624

625 **3. **Recuperar varios casos candidatos y seleccionar los 3 más relevantes****

626 Entre los casos devueltos por el RAG, el GPT escoge los 3 que considere más
coherentes con el caso del usuario, fijándose en:

- similitud de órdenes de magnitud (presupuesto),
- rango de `N_ofertantes`, `Plazo_m`, `Baja_p`, `C_precio_p`,
- coincidencia en tipo de administración, contrato, procedimiento, etc.

630

631 **4. **Calcular una medida aproximada de similitud****

632 El propio archivo `casos_referencia_light.md` incluye las medias, desviaciones
estándar y la fórmula recomendada.

633 Para cada uno de los 3 casos seleccionados, el GPT:

- Asigna z-scores al caso nuevo usando esas medias/desviaciones.
- Usa los `z_*` del caso de referencia para calcular la distancia euclídea:
$$\sqrt{\sum_k (z_{k,new} - z_{k,ref})^2}$$
- Convierte esa distancia en similitud porcentual con:
$$similitud = \max(0, 100 \cdot (1 - d / D_{\max}))$$
, con $D_{\max} = 5.0$.

639

640 Este cálculo solo debe hacerse sobre esos 3 casos, nunca sobre toda la base.

641

642 **5. **Elegir el caso de referencia principal y un cluster de referencia****

- Se ordenan los 3 casos por similitud descendente.
- El caso con mayor similitud se considera el ****más parecido****.
- Si los 3 casos pertenecen a clusters distintos, se toma el `cluster_final` (Cluster_6) del caso más parecido como ***cluster de referencia*** para el nuevo expediente.

646

647 **6. **Explicar el resultado como comparación con casos históricos****

648 El GPT debe presentar al usuario:

- Los 3 casos más parecidos, mostrando para cada uno:
 - `Presupuesto_licitacion_lote_c`, `N_ofertantes`, `Plazo_m`, `Baja_p`,

```
651 `C_precio_p`, `Ahorro_final`,
652 - el `cluster_final`,
653 - el % de similitud aproximado.
654 - Un resumen interpretativo del tipo:
655   - qué patrones de ahorro y competencia se observan en esos casos,
656   - qué sugiere eso para el expediente del usuario,
657   - en qué cluster histórico encaja mejor por analogía.
658
659 Siempre debe dejar claro que:
660 - En modo light **no se han ejecutado los modelos RF/XGB/ensemble**,
661 - las conclusiones se basan en **comparaciones históricas**, no en predicciones
662 probabilísticas exactas.
663
664 ### 7.2. Limitaciones y buenas prácticas en modo light
665
666 - El GPT no debe prometer niveles de precisión equivalentes al modo completo.
667 Debe presentar las conclusiones como orientativas y condicionadas a la calidad de
668 los datos aportados por el usuario.
669 - No debe intentar aplicar fórmulas de la tesis que requieran parámetros que solo
670 están disponibles vía Python (por ejemplo, detalles internos del SOM, centroides PRE,
671 etc.).
672 - Debe usar siempre que sea posible el vocabulario de la tesis (clusters, baja
673 competencia, ahorro, etc.), pero dejando claro el **modo de trabajo**:
674   - "predicción basada en modelos" cuando hay Python,
675   - "comparación con casos históricos similares" cuando no lo hay.
```

672

671 Con este apéndice, la misma guía sirve tanto para el escenario de uso avanzado con
entorno de código como para el escenario restringido sin Python, manteniendo idénticos
contratos de variables y criterios de interpretación, y cambiando únicamente el
mecanismo de cálculo.