

```

1  ### SCRIPT 12 - R ###
2  #
3  # Análisis SOM no supervisado (5 variables) + K-means (K = 8)
4  # Mapeo de variable objetivo Ahorro_final = (Presupuesto - Precio_final)/Presupuesto
5  #
6  # =====
7  #
8  # ---- Paquetes -----
9
10 dep_pkgs <- c(
11   "kohonen",
12   "dplyr",
13   "tidyverse",
14   "readxl",
15   "readr",
16   "ggplot2",
17   "tibble",
18   "stringr",
19   "colorspace"
20 )
21
22 invisible(lapply(dep_pkgs, function(p) {
23   if (!requireNamespace(p, quietly = TRUE)) install.packages(p)
24 }))
25
26 lapply(dep_pkgs, library, character.only = TRUE) |> invisible()
27
28 set.seed(321)
29 options(dplyr.summarise.inform = FALSE)
30
31 # ---- Parámetros y rutas -----
32
33 root_dir <- "F:/Tesis Guillermo_casa/Tesis Guillermo/_BD actualizada/Datos
atom/PLACSP/_final_Modelos"
34 file_train <- file.path(root_dir, "df_train_v5_completo.xlsx")
35 file_test <- file.path(root_dir, "df_test_v5_completo.xlsx")
36
37 # Variables exactas del dataset
38 vars_som <- c(
39   "Presupuesto_licitacion_lote_c",
40   "N_ofertantes",
41   "C_precio_p",
42   "Plazo_m",
43   "Baja_p"
44 )
45
46 # Objetivo a mapear en el SOM
47 obj_label <- "Ahorro_final" # (Presupuesto - Precio_final) / Presupuesto
48 req_obj_cols <- c("Presupuesto_licitacion_lote_c", "Precio_final_e_c")
49
50 # Variables disponibles en licitación para XYF
51 vars_licit <- c(
52   "N_lotes",
53   "Intervalo_urge_d",
54   "N_clasi_empresa",
55   "N_CPV",
56   "Tipo_de_Administracion_c",
57   "Tipo_de_procedimiento_c"
58 )
59
60 k_som <- 8
61 som_xy <- c(10, 10)
62
63 # ---- Utilidades -----
64
65 check_required_cols <- function(df, cols, name) {
66   miss <- setdiff(cols, names(df))
67   if (length(miss) > 0) stop(sprintf("Faltan columnas en %s: %s", name, paste(miss,
68   collapse = ", ")))
69 }
70
71 scale_fit <- function(x) {
72   m <- scale(x)

```

```

72     list(x = m, params = list(center = attr(m, "scaled:center"), scale = attr(m,
73         "scaled:scale"))))
74   }
75 
76   scale_apply <- function(x, params) {
77     scale(x, center = params$center, scale = params$scale)
78   }
79 
80   make_mm_licit <- function(df) {
81     df2 <- df |>
82       mutate(
83         Tipo_de_Administracion_c = as.factor(Tipo_de_Administracion_c),
84         Tipo_de_procedimiento_c = as.factor(Tipo_de_procedimiento_c)
85       )
86     mm <- stats::model.matrix(
87       ~ 0 + N_lotes + Intervalo_urgen_d + N_clasi_empresa + N_CPV +
88       Tipo_de_Administracion_c + Tipo_de_procedimiento_c,
89       data = df2
90     )
91     storage.mode(mm) <- "double"
92     mm
93   }
94 
95   align_cols <- function(x_new, col_ref) {
96     miss <- setdiff(col_ref, colnames(x_new))
97     if (length(miss) > 0) x_new <- cbind(x_new, matrix(0, nrow = nrow(x_new), ncol =
98       length(miss), dimnames = list(NULL, miss)))
99     extra <- setdiff(colnames(x_new), col_ref)
100    if (length(extra) > 0) x_new <- x_new[, setdiff(colnames(x_new), extra), drop =
101      FALSE]
102    x_new[, col_ref, drop = FALSE]
103  }
104 
105  majority_label <- function(x) { x <- as.character(x); tab <- table(x); names(tab)[
106    which.max(tab)] }
107 
108  # ----- Carga de datos -----
109 
110  message("Leyendo train/test desde: ", root_dir)
111 
112  df_train <- readxl::read_xlsx(file_train) |> tibble::as_tibble()
113  df_test <- readxl::read_xlsx(file_test) |> tibble::as_tibble()
114 
115  check_required_cols(df_train, c(vars_som, req_obj_cols, vars_licit), "train")
116  check_required_cols(df_test, c(vars_som, req_obj_cols, vars_licit), "test")
117 
118  # ----- Preparación para SOM (5 variables) -----
119 
120  x_train_raw <- df_train |>
121    select(all_of(vars_som)) |>
122    mutate(across(everything(), as.numeric))
123 
124  x_test_raw <- df_test |>
125    select(all_of(vars_som)) |>
126    mutate(across(everything(), as.numeric))
127 
128  keep_train <- stats::complete.cases(x_train_raw)
129  keep_test <- stats::complete.cases(x_test_raw)
130 
131  message(sprintf("SOM | Train completas: %d/%d | Test completas: %d/%d",
132    sum(keep_train), nrow(df_train), sum(keep_test), nrow(df_test)))
133 
134  x_train_raw <- x_train_raw[keep_train, , drop = FALSE]
135  x_test_raw <- x_test_raw[keep_test, , drop = FALSE]
136 
137  # Objetivo = Ahorro_final
138  pli <- df_train$Presupuesto_licitacion_lote_c[keep_train]
139  pfin <- df_train$Precio_final_e_c[keep_train]
140  y_train_obj <- (pli - pfin) / pli

```

```

139 y_train_obj[!is.finite(y_train_obj)] <- NA_real_
140
141 # Escalado SOM
142 sc_som <- scale_fit(as.matrix(x_train_raw))
143 x_train <- sc_som$x
144 x_test <- scale_apply(as.matrix(x_test_raw), sc_som$params)
145
146 # ---- Entrenamiento SOM -----
147
148 som_grid <- kohonen::somgrid(xdim = som_xy[1], ydim = som_xy[2], topo = "hexagonal",
149 toroidal = FALSE)
150
151 set.seed(321)
152 som_model <- kohonen::som(
153   X = x_train,
154   grid = som_grid,
155   rlen = 10000,
156   alpha = c(0.05, 0.005),
157   keep.data = TRUE
158 )
159
160 # K-means sobre códigos del SOM
161 set.seed(321)
162 codes_matrix <- som_model$codes[[1]]
163 km_codes <- stats::kmeans(codes_matrix, centers = k_som, nstart = 50, iter.max = 100)
164 unit_to_cluster <- km_codes$cluster
165
166 eur_train <- som_model$unit.classif
167 cl_train_som <- unit_to_cluster[eur_train]
168
169 # ---- Resúmenes (neurona/cluster) -----
170
171 obj_por_neur <- tapply(y_train_obj, eur_train, mean, na.rm = TRUE)
172 prop_vec <- rep(NA_real_, som_xy[1] * som_xy[2])
173 prop_vec[as.integer(names(obj_por_neur))] <- as.numeric(obj_por_neur)
174
175 cluster_mean_obj <- tibble(
176   cluster = factor(sort(unique(cl_train_som))),
177   obj_medio = as.numeric(tapply(y_train_obj, cl_train_som, mean, na.rm = TRUE)[levels(
178     factor(cl_train_som))])
179 )
180
181 resumen_clusters <- tibble(cluster = factor(cl_train_som), obj = y_train_obj) |>
182   dplyr::group_by(cluster) |>
183   dplyr::summarise(
184     n_casos = dplyr::n(),
185     obj_media = mean(obj, na.rm = TRUE),
186     obj_sd = sd(obj, na.rm = TRUE)
187   ) |>
188   dplyr::ungroup() |>
189   arrange(as.integer(as.character(cluster)))
190
191 # ---- Visualización -----
192
193 pdf(file.path(root_dir, "SOM_K8_Simple_Mapas.pdf"), width = 8.3, height = 11.7)
194 par(mfrow = c(4, 4), mar = c(2, 2, 3, 1))
195
196 plot(som_model, type = "dist.neighbours", main = "U-Matrix", palette.name = pal_cf,
197 shape = "straight")
198 add.cluster.boundaries(som_model, unit_to_cluster, lwd = 1.75, col = "white")
199
200 for (i in seq_len(ncol(x_train))) {
201   plot(som_model, type = "property", property = som_model$codes[[1]][, i],
202         main = colnames(x_train)[i], palette.name = pal_cf, shape = "straight")
203   add.cluster.boundaries(som_model, unit_to_cluster, lwd = 1.75, col = "white")
204 }
205
206 plot(som_model, type = "count", main = "Casos por neurona", palette.name = pal_cf,
207 shape = "straight")
208 add.cluster.boundaries(som_model, unit_to_cluster, lwd = 1.75, col = "white")
209

```

```

207     main = paste0("Media ", obj_label, " por neurona (train)"),
208     palette.name = pal_cf, shape = "straight")
209 add.cluster.boundaries(som_model, unit_to_cluster, lwd = 1.75, col = "white")
210
211 dev.off()
212
213 # Pie chart de componentes
214 pdf(file.path(root_dir, "SOM_K8_Simple_Codes_PieChart.pdf"), width = 8.3, height =
215 11.7)
215 plot(som_model, type = "codes", main = "SOM K8 - Componentes (Pie Chart)",
216 palette.name = hi_contrast, shape = "straight")
216 add.cluster.boundaries(som_model, unit_to_cluster, lwd = 1.75, col = "white")
217 dev.off()
218
219 # Fondo con objetivo
220 pdf(file.path(root_dir, "SOM_K8_Simple_Background_AhorroFinal.pdf"), width = 8.3,
220 height = 11.7)
221 plot(som_model, type = "property", property = prop_vec, main = paste0("Media ",
222 obj_label, " por neurona (train)"), palette.name = pal_cf, shape = "straight")
222 add.cluster.boundaries(som_model, unit_to_cluster, lwd = 1.75, col = "white")
223 dev.off()
224
225 # ---- Predicción con SOM -----
226
227 predict_som_cluster <- function(new_x_scaled) {
228   if (is.null(dim(new_x_scaled)) || any(dim(new_x_scaled) == 0)) {
229     return(tibble(neurona_bmu = integer(0), cluster_som = factor(character(0), levels
229      = sort(unique(unit_to_cluster)))))
230   }
231   mapped <- predict(som_model, newdata = as.matrix(new_x_scaled))
232   bmu <- as.integer(mapped$unit.classif)
233   cl <- unit_to_cluster[bmu]
234   tibble(neurona_bmu = bmu, cluster_som = factor(cl, levels = sort(unique(
234     unit_to_cluster))))
235 }
236
237 add_expected_obj <- function(tbl_clusters) {
238   tbl_clusters |>
239     left_join(cluster_mean_obj, by = c("cluster_som" = "cluster")) |>
240     rename(esperado_cluster = obj_medio)
241 }
242
243 # ---- Clasificación con XYF (predict de Kohonen) -----
244
245 train_som_rows <- df_train[keep_train, , drop = FALSE]
246 keep_xyf_train <- stats::complete.cases(train_som_rows[, vars_licit])
247 message(sprintf("XYF | Train completas (licitación): %d/%d", sum(keep_xyf_train), sum(
247   keep_train)))
248
249 if (sum(keep_xyf_train) > 0) {
250   y_train_xyf <- factor(cl_train_som[keep_xyf_train])
251
252   x_licit_train_mm <- make_mm_licit(train_som_rows[keep_xyf_train, vars_licit, drop =
252     FALSE])
253   sc_xyf <- scale_fit(x_licit_train_mm)
254   x_licit_train <- sc_xyf$x
255   mm_cols <- colnames(x_licit_train)
256
257   xyf_grid <- kohonen::somgrid(xdim = som_xy[1], ydim = som_xy[2], topo = "hexagonal"
257     , toroidal = FALSE)
258
259   set.seed(321)
260   xyf_model <- kohonen::xyf(
261     X = x_licit_train,
262     Y = y_train_xyf,
263     grid = xyf_grid,
264     rlen = 5000,
265     alpha = c(0.05, 0.005),
266     keep.data = TRUE
267   )
268
269   unit_xyf_train <- xyf_model$unit.classif
270   unit_labels <- tibble(unit = unit_xyf_train, label = y_train_xyf) |>

```

```

271 summarise(label = majority_label(label), .by = unit)
272
273 } else {
274   message("Aviso: no hay filas completas en TRAIN para variables de licitación; no se
275   entrena XYF.")
276   xyf_model <- NULL
277   sc_xyf <- list(params = list(center = NULL, scale = NULL))
278   x_licit_train <- matrix(numeric(0), nrow = 0, ncol = 0)
279   mm_cols <- character(0)
280   xyf_grid <- NULL
281   unit_labels <- tibble(unit = integer(0), label = factor(character(0)))
282 }
283
283 predict_xyf_cluster <- function(new_mm_scaled) {
284   # Devuelve tantas filas como new_mm_scaled tenga; NAs si no se puede predecir
285   if (is.null(xfy_model)) {
286     n_rows <- if (!is.null(new_mm_scaled)) NROW(new_mm_scaled) else 0L
287     return(tibble(
288       neurona_bmu_xyf = rep(NA_integer_, n_rows),
289       cluster_xyf = factor(rep(NA_character_, n_rows), levels = levels(y_train_xyf))
290     ))
291   }
292   if (is.null(new_mm_scaled)) {
293     return(tibble(neurona_bmu_xyf = integer(0),
294                   cluster_xyf = factor(character(0), levels = levels(y_train_xyf))))
295   }
296   new_mm_scaled <- as.matrix(new_mm_scaled)
297   n_rows <- nrow(new_mm_scaled)
298   if (is.null(dim(new_mm_scaled)) || n_rows == 0 || ncol(new_mm_scaled) == 0) {
299     return(tibble(neurona_bmu_xyf = integer(0),
300                   cluster_xyf = factor(character(0), levels = levels(y_train_xyf))))
301   }
302   storage.mode(new_mm_scaled) <- "double"
303   pr <- tryCatch(
304     predict(xfy_model, newdata = list(X = new_mm_scaled), whatmap = 1, maxNA.fraction
305     = 1),
306     error = function(e) { message("Aviso: predict(XYF) omitido: ", e$message); NULL }
307   )
308   if (is.null(pr)) {
309     return(tibble(
310       neurona_bmu_xyf = rep(NA_integer_, n_rows),
311       cluster_xyf = factor(rep(NA_character_, n_rows), levels = levels(y_train_xyf))
312     ))
313   }
314   bmu <- as.integer(pr$unit.classif)
315   lab <- unit_labels$label[match(bmu, unit_labels$unit)]
316   tibble(neurona_bmu_xyf = bmu, cluster_xyf = factor(lab, levels = levels(y_train_xyf)))
317 }
318
318 # ----- Predicciones en TEST -----
319
320 if (sum(keep_test) > 0) {
321   som_test_rows <- tibble(.row_id = which(keep_test)) %>
322     bind_cols(predict_som_cluster(x_test)) %>
323     add_expected_obj()
324 } else {
325   som_test_rows <- tibble(.row_id = integer(0), neurona_bmu = integer(0),
326                           cluster_som = factor(character(0), levels = sort(unique(
327                             unit_to_cluster))),
328                           esperado_cluster = numeric(0))
329 }
330
330 keep_xyf_test <- stats:::complete.cases(df_test[, vars_licit])
331 message(sprintf("XYF | Test completas (licitación): %d/%d", sum(keep_xyf_test), nrow(
332   df_test)))
332 if (sum(keep_xyf_test) > 0 && !is.null(xfy_model)) {
333   x_licit_test_mm_full <- make_mm_licit(df_test[keep_xyf_test, vars_licit, drop =
334     FALSE])
335   # Alinear columnas del test a las del train (mm_cols)
336   x_licit_test_mm <- align_cols(x_licit_test_mm_full, mm_cols)
337   # Debug opcional: categorías faltantes o extra

```

```

338 miss_cols <- setdiff(mm_cols, colnames(x_licit_test_mm_full))
339 extra_cols <- setdiff(colnames(x_licit_test_mm_full), mm_cols)
340 if (length(miss_cols) > 0) message("XYF | Dummies faltantes en TEST añadidas como
0: ", paste(miss_cols, collapse = ", "))
341 if (length(extra_cols) > 0) message("XYF | Dummies extra en TEST ignoradas: ", paste
(extra_cols, collapse = ", "))
342
343 if (nrow(x_licit_test_mm) > 0 && ncol(x_licit_test_mm) > 0) {
344   x_licit_test <- scale_apply(as.matrix(x_licit_test_mm), sc_xyf$params)
345
346   xyf_preds <- predict_xyf_cluster(x_licit_test)
347   xyf_test_rows <- tibble(.row_id = which(keep_xyf_test)) %>
348     bind_cols(xyf_preds)
349 } else {
350   message("Aviso: matriz XYF de TEST sin dimensión válida; se omiten predicciones
XYF.")
351   xyf_test_rows <- tibble(.row_id = integer(0), neurona_bmu_xyf = integer(0),
352                           cluster_xyf = factor(character(0), levels = sort(unique(
353                             cl_train_som))))
354 }
355 } else {
356   if (sum(keep_xyf_test) == 0) message("Aviso: ninguna fila de TEST completa para
licitación; se omiten predicciones XYF.")
357   xyf_test_rows <- tibble(.row_id = integer(0), neurona_bmu_xyf = integer(0),
358                           cluster_xyf = factor(character(0), levels = sort(unique(
359                             cl_train_som))))
360 }
361
362 pred_test <- tibble(.row_id = seq_len(nrow(df_test))) %>
363   left_join(som_test_rows, by = ".row_id") %>
364   left_join(xyf_test_rows, by = ".row_id")
365
366 # ----- Evaluación en TEST: Confusión y Scatter -----
367 # Ahorro real
368 ahorro_real_test <- with(df_test, (Presupuesto_licitacion_lote_c - Precio_final_e_c) /
369   Presupuesto_licitacion_lote_c)
370 ahorro_real_test[!is.finite(ahorro_real_test)] <- NA_real_
371
372 # Añadir real y renombrar predicción de SOM
373 pred_test <- pred_test %>
374   mutate(
375     ahorro_real = ahorro_real_test,
376     ahorro_pred_som = esperado_cluster
377   )
378
379 # Predicción de ahorro por XYF (media del cluster XYF)
380 pred_test <- pred_test %>
381   left_join(cluster_mean_obj %>
382     rename(cluster_xyf = cluster, ahorro_pred_xyf = obj_medio),
383     by = "cluster_xyf")
384
385 # Matriz de confusión SOM vs XYF (solo filas con ambos clusters)
386 idx_cm <- stats::complete.cases(pred_test$cluster_som, pred_test$cluster_xyf)
387 if (any(idx_cm)) {
388   cm <- table(SOM = pred_test$cluster_som[idx_cm], XYF = pred_test$cluster_xyf[idx_cm])
389   acc <- mean(pred_test$cluster_som[idx_cm] == pred_test$cluster_xyf[idx_cm])
390   cm_df <- as.data.frame.matrix(cm)
391   readr::write_csv(tibble::rownames_to_column(as.data.frame(cm_df), var = "SOM"),
392     file.path(root_dir, "SOM_K8_Simple_confusion_SOM_vs_XYF.csv"))
393 } else {
394   acc <- NA_real_
395 }
396
397 # Métricas de ajuste (real vs predicho)
398 mae <- function(a, p) { idx <- stats::complete.cases(a, p); if (!any(idx)) return(
399   NA_real_); mean(abs(a[idx] - p[idx])) }
400 rmse <- function(a, p) { idx <- stats::complete.cases(a, p); if (!any(idx)) return(
401   NA_real_); sqrt(mean((a[idx] - p[idx])^2)) }
402 correl <- function(a, p) { idx <- stats::complete.cases(a, p); if (sum(idx) < 2)
403   return(NA_real_); suppressWarnings(cor(a[idx], p[idx])) }
404
405 metrics_test <- dplyr::bind_rows(

```

```

400 tibble(modelo = "SOM",
401     MAE = mae(pred_test$ahorro_real, pred_test$ahorro_pred_som),
402     RMSE = rmse(pred_test$ahorro_real, pred_test$ahorro_pred_som),
403     COR = correl(pred_test$ahorro_real, pred_test$ahorro_pred_som)),
404 tibble(modelo = "XYF",
405     MAE = mae(pred_test$ahorro_real, pred_test$ahorro_pred_xyf),
406     RMSE = rmse(pred_test$ahorro_real, pred_test$ahorro_pred_xyf),
407     COR = correl(pred_test$ahorro_real, pred_test$ahorro_pred_xyf))
408 ) |>
409     mutate(ACC_clusters_SOM_vs_XYF = acc)
410
411 readr::write_csv(metrics_test, file.path(root_dir, "SOM_K8_Simple_test_metrics.csv"))
412
413 # Log dimensiones para depuración
414 message(sprintf("XYF dims | train: %d x %d | test: %d x %d | preds: %d filas",
415                 ncol(x_licit_train), nrow(x_licit_train),
416                 if (exists("x_licit_test")) ncol(x_licit_test) else -1,
417                 if (exists("x_licit_test")) nrow(x_licit_test) else -1,
418                 nrow(xyf_test_rows)))
419
420 # Scatter real vs predicho
421 plot_df <- tidyverse::pivot_longer(
422     pred_test,
423     cols = c(ahorro_pred_som, ahorro_pred_xyf),
424     names_to = "metodo",
425     values_to = "ahorro_pred"
426 )
427
428 p_sc <- ggplot(plot_df, aes(x = ahorro_real, y = ahorro_pred, color = metodo)) +
429     geom_point(alpha = 0.6) +
430     geom_abline(slope = 1, intercept = 0, linetype = "dashed", color = "grey40") +
431     scale_color_manual(values = c(ahorro_pred_som = "steelblue", ahorro_pred_xyf =
432         "darkorange"),
433                         labels = c(ahorro_pred_som = "SOM", ahorro_pred_xyf = "XYF")) +
434     labs(title = "TEST | Ahorro real vs predicho", x = "Ahorro real", y = "Ahorro
435     predicho", color = "Modelo") +
436     theme_minimal()
437
438 ggsave(file.path(root_dir, "SOM_K8_Simple_Test_Scatter_Ahorro.pdf"), plot = p_sc,
439         width = 8, height = 6)
440
441 # ---- Guardado -----
442
443 saveRDS(list(
444     som_model      = som_model,
445     som_grid       = som_grid,
446     kmeans_codes   = km_codes,
447     unit_to_cluster = unit_to_cluster,
448     scaler_params_som= sc_som$params,
449     vars_som        = vars_som,
450     cluster_mean_obj = cluster_mean_obj,
451     objetivo_label = obj_label
452 ), file.path(root_dir, "SOM_K8_Simple_model.rds"))
453
454 saveRDS(list(
455     xyf_model      = xyf_model,
456     xyf_grid       = xyf_grid,
457     scaler_params_xyf= sc_xyf$params,
458     mm_cols_xyf    = mm_cols,
459     unit_labels_xyf = unit_labels,
460     vars_llicit    = vars_llicit
461 ), file.path(root_dir, "SOM_K8_Simple_xyf_model.rds"))
462
463 readr::write_csv(resumen_clusters, file.path(root_dir,
464 "SOM_K8_Simple_resumen_clusters.csv"))
465 readr::write_csv(pred_test,           file.path(root_dir, "SOM_K8_Simple_pred_TEST.csv"))
466
467 cat("\n== Resumen por cluster (SOM + K=8) ==\n"); print(resumen_clusters)
468 cat("\nArchivos guardados en: ", root_dir, "\n")
469
470 # ----- FIN -----

```