

```

1  ##### SCRIPT 5 - R #####
2  # Obtención de plazos de ejecucion desde ATOM
3  #
4  ##########
5  library(XML)
6  library(dplyr)
7  library(foreach)
8  library(doParallel)
9
10 # Configuración paralela
11 cores <- detectCores()
12 cl <- makeCluster(cores[1] - 1, type = "PSOCK", outfile = "log.txt")
13 registerDoParallel(cl)
14
15 # Ruta de trabajo
16 setwd("C:/Users/guillermo.alonso/Desktop/Tesis Guillermo/_BD actualizada/Datos atom/PLACSP")
17 PATH <- "licitacionesPerfilesContratanteCompleto3_2025" # carpeta donde están los .atom
18 RData_FILE <- "3.34_PlannedPeriod_Extension_2025.Rdata"
19 DBtablename <- "PlannedPeriodExtension"
20
21 # Función para encontrar un nodo padre por nombre
22 DamePadre <- function(xml_a, tag) {
23   xml_p <- xml_a
24   repeat {
25     xml_p <- xmlParent(xml_p)
26     if (xmlName(xml_p) == tag) break
27   }
28   return(xml_p)
29 }
30
31 # Función segura para extraer el valor de un nodo XML
32 safe_xml_value <- function(node, tag) {
33   if (!is.null(node) && !is.null(node[[tag]])) {
34     return(xmlValue(node[[tag]]))
35   } else {
36     return(NA)
37   }
38 }
39
40 # Columnas objetivo
41 nombres <- c("ContractID", "entryID",
42             "start_date", "end_date",
43             "duration", "duration_unit",
44             "description", "extension_description",
45             "updated")
46
47 # Inicializar data.frame vacío
48 Proy <- data.frame(matrix(NA, 0, length(nombres)))
49 colnames(Proy) <- nombres
50
51 # Listar archivos atom
52 arch <- list.files(path = PATH, pattern = "\\.atom$", full.names = FALSE)
53
54 # Procesar en paralelo
55 Proy <- foreach(ci = 1:length(arch), .combine = rbind, .packages = c("XML", "dplyr"))
%dopar%
56
57   archivo_temp <- xmlParse(file.path(PATH, arch[ci]))
58   cat(sprintf("[%s] Procesando archivo %s\n", Sys.time(), arch[ci]))
59
60   buscar <- xmlElementsByTagName(xmlRoot(archivo_temp), "ProcurementProject",
61                                   recursive = TRUE)
62
63   Proy <- data.frame(matrix(NA, 0, length(nombres)))
64   colnames(Proy) <- nombres
65
66   for (xmla in buscar) {
67     foo <- Proy[1, ]
68
69     tryCatch({
      xmfp <- DamePadre(xmla, "ContractFolderStatus")

```

```

70 foo$ContractID <- safe_xml_value(xmlp, "ContractFolderID")
71
72 entry_node <- xmlParent(xmlp)
73 id_value <- safe_xml_value(entry_node, "id")
74 updated_raw <- safe_xml_value(entry_node, "updated")
75
76 if (!is.na(id_value)) {
77   bar <- strsplit(id_value, "/")[[1]]
78   foo$entryID <- bar[length(bar)]
79 }
80
81 foo$updated <- tryCatch(as.Date(updated_raw), error = function(e) NA)
82
83 # PlannedPeriod
84 planned_period <- xmla[["PlannedPeriod"]]
85 if (!is.null(planned_period)) {
86   foo$start_date <- safe_xml_value(planned_period, "StartDate")
87   foo$end_date <- safe_xml_value(planned_period, "EndDate")
88   foo$duration <- safe_xml_value(planned_period, "DurationMeasure")
89   foo$duration_unit <- if (!is.null(planned_period[["DurationMeasure"]])) {
90     xmlGetAttr(planned_period[["DurationMeasure"]], "unitCode")
91   } else {
92     NA
93   }
94   foo$description <- safe_xml_value(planned_period, "Description")
95 }
96
97 # ContractExtension
98 contract_ext <- xmla[["ContractExtension"]]
99 if (!is.null(contract_ext)) {
100   option_period <- contract_ext[["OptionValidityPeriod"]]
101   foo$extension_description <- safe_xml_value(option_period, "Description")
102 }
103
104 Proy <- rbind(Proy, foo)
105 }, error = function(e) {
106   cat(sprintf("⚠ Error en archivo %s: %s\n", arch[ci], e$message))
107 }
108 }
109
110 Proy
111 }
112
113 # Finalizar clúster
114 stopCluster(cl)
115
116 ### Limpieza y tratamiento de resultados en el DF Proy -----
117
118 ## Eliminamos casos vacíos
119 Proy <- Proy %>%
120   filter(!is.na(start_date) & is.na(duration))
121
122 ## Damos formato a las variables
123
124 Proy <- Proy |>
125   mutate(
126     start_date = as.Date(start_date),
127     end_date = as.Date(end_date),
128     updated = as.Date(updated),
129     duration = suppressWarnings(as.numeric(duration)),
130     description = na_if(trimws(description), ""),
131     extension_description = na_if(trimws(extension_description), ""))
132 )
133 ## Calculamos variable "duration" donde no la hay en formato MES (MON)
134 Proy <- Proy %>%
135   mutate(
136     duration = ifelse(
137       is.na(duration) & !is.na(start_date) & !is.na(end_date),
138       floor(as.numeric(end_date - start_date) / 30),
139       duration
140     ),
141     duration = ifelse(duration == 0, 1, duration),
142     duration = as.integer(duration),

```

```

143     duration_unit = ifelse(
144         is.na(duration_unit) & !is.na(start_date) & !is.na(end_date),
145         "MON",
146         duration_unit
147     )
148   )
149 
150 ## Pasamos todas las duraciones a meses (MON) - 30 dias y 12 meses
151 Proy <- Proy %>%
152   mutate(
153     duration_m = case_when(
154       duration_unit == "DAY" ~ floor(duration / 30),
155       duration_unit == "ANN" ~ duration * 12,
156       duration_unit == "MON" ~ duration,
157       TRUE ~ NA_integer_
158     ),
159     duration_m = ifelse(duration_m == 0 & !is.na(duration), 1L, duration_m),
160     duration_m = as.integer(duration_m)
161   )
162 
163 ## Selecciono variables interesantes
164 library(stringr)
165 
166 Proy <- Proy %>%
167   transmute(
168     Identificador = entryID,
169     Identificador_lici = str_c(entryID, coalesce(as.character(ContractID), "NA"), sep
170     = "//"),
171     Fecha_actualizacion = updated,
172     Plazo_m = duration_m,
173     Prorroga = extension_description
174   )
175 
176 ## Hacemos un unique
177 Proy <- unique(Proy)
178 
179 ## En los casos que Identificador_lici sea el mismo, nos quedamos con el caso de
180 ## Fecha_actualizacion mas tardia
181 Proy <- Proy %>%
182   group_by(Identificador_lici) %>%
183   slice_max(order_by = Fecha_actualizacion, n = 1, with_ties = FALSE) %>%
184   ungroup()
185 
186 # Guardar resultados
187 save(Proy, nombres, file = RData_FILE).
188 
189 library(ggplot2)
190 library(dplyr)
191 library(openxlsx)
192 library(RODBC)
193 library(quanteda) #paquete recomendado en todos los
194 library(quanteda.textmodels) #aux
195 library(quanteda.textstats) # aux
196 library(readtext) # sencilla manera de leer data de texto en R
197 library(spacyr) # NLP usando la libreria spaCy, incluyendo etiquetado part-of-speech,
198 entity recognition y dependency parsing.
199 library(devtools)
200 library(zoo)
201 library(stringr)
202 # Ruta de trabajo
203 setwd("C:/Users/guillermo.alonso/Desktop/Tesis Guillermo/_BD actualizada/Datos
204 atom/PLACSP")
205 ##### Generamos data frame de Plazo de proyecto DF_plazo -----
206 # Carga de archivos por año
207 
208 años <- 2012:2025
209 lista_df <- list()
210 
211 for (año in años) {
  archivo <- paste0("3.34_PlannedPeriod_Extension_", año, ".RData")
  load(archivo)
  lista_df[[as.character(año)]] <- Proy

```

```
212 }
213
214 # Unificación de datos en DF_plazo_u
215 DF_plazo <- bind_rows(lista_df)
216 DF_plazo_u <- unique(DF_plazo) # Eliminamos duplicados
217 DF_plazo_u <- DF_plazo_u %>%
218   group_by(Identificador_lici) %>%
219   slice_max(order_by = Fecha_actualizacion, n = 1, with_ties = FALSE) %>%
220   ungroup()
221
222 ## Guardo resultados de plazo
223
224 save(DF_plazo_u, file="DF_plazo_u.RData")
225 write.xlsx(DF_plazo_u,"DF_plazo_u.xlsx", colNames = T, rowNames=F)
226
227
```