

```

1  ### SCRIPT 13 - PYTHON
2  # ANÁLISIS CLUSTERING COMPARATIVA K=6 y K=8
3  #
4  # =====
5  # IMPORTACIÓN DE LIBRERÍAS
6  # =====
7
8  import pandas as pd
9  import numpy as np
10 from datetime import datetime
11 from sklearn.cluster import KMeans
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.decomposition import PCA
14 from sklearn.metrics import silhouette_score
15 import warnings
16 warnings.filterwarnings('ignore')
17
18 print("*"*80)
19 print("PROYECTO: ANÁLISIS CLUSTERING LICITACIONES PÚBLICAS V8")
20 print("*"*80)
21 print("Dataset combinado Train+Test con clustering K=6 y K=8")
22
23 # =====
24 # PASO 1: CARGA DE DATASETS BRUTOS
25 # =====
26
27 print("\n1. CARGA DE DATASETS BRUTOS:")
28 print("-" * 50)
29
30 # Cargar datasets originales
31 df_train = pd.read_excel('df_train_v5_completo.xlsx')
32 df_test = pd.read_excel('df_test_v5_completo.xlsx')
33
34 print(f"\nDataset entrenamiento cargado: {df_train.shape}")
35 print(f"\nDataset test cargado: {df_test.shape}")
36 print(f"\nTotal registros: {df_train.shape[0] + df_test.shape[0]}")
37
38 # Verificar columnas clave
39 columnas_clave = ['Primera_publicacion_c', 'Presupuesto_licitacion_lote_c',
40                   'Precio_final_e_c', 'C_precio_p']
41 print(f"\nVerificación columnas clave:")
42 for col in columnas_clave:
43     train_exists = col in df_train.columns
44     test_exists = col in df_test.columns
45     print(f" {col}: Train={train_exists}, Test={test_exists}")
46
47 # =====
48 # PASO 2: CREACIÓN DE VARIABLES DERIVADAS
49 # =====
50
51 print(f"\n2. CREACIÓN DE VARIABLES DERIVADAS:")
52 print("-" * 50)
53
54 # Variable Ahorro_final con fórmula correcta
55 # Ahorro_final (%) = (Presupuesto_licitacion_lote_c - Precio_final_e_c) / Presupuesto_licitacion_lote_c * 100
56 df_train['Ahorro_final'] = ((df_train['Presupuesto_licitacion_lote_c'] - df_train['Precio_final_e_c']) /
57                               df_train['Presupuesto_licitacion_lote_c'] * 100)
58
59 df_test['Ahorro_final'] = ((df_test['Presupuesto_licitacion_lote_c'] - df_test['Precio_final_e_c']) /
60                               df_test['Presupuesto_licitacion_lote_c'] * 100)
61
62 # Variable Mes_lici desde Primera_publicacion_c
63 df_train['Mes_lici'] = df_train['Primera_publicacion_c'].dt.month
64 df_test['Mes_lici'] = df_test['Primera_publicacion_c'].dt.month
65
66 print(f"\nVariables derivadas creadas exitosamente")
67
68 # Estadísticas de Ahorro_final
69 print(f"\nESTADÍSTICAS AHORRO_FINAL:")

```

```

70 print(f"Train - Media: {df_train['Ahorro_final'].mean():.2f}%)")
71 print(f"Train - Desv.Std: {df_train['Ahorro_final'].std():.2f}%)")
72 print(f"Train - Sobrecostos: {(df_train['Ahorro_final'] < 0).sum()} ({(df_train['Ahorro_final'] < 0).mean()*100:.1f}%)")
73
74 print(f"Test - Media: {df_test['Ahorro_final'].mean():.2f}%)")
75 print(f"Test - Desv.Std: {df_test['Ahorro_final'].std():.2f}%)")
76 print(f"Test - Sobrecostos: {(df_test['Ahorro_final'] < 0).sum()} ({(df_test['Ahorro_final'] < 0).mean()*100:.1f}%)")
77 # =====
78 # PASO 3: UNIÓN DE DATASETS
79 # =====
80
81 print(f"\n3. UNIÓN DE DATASETS:")
82 print("-" * 50)
83
84
85 # Crear variable identificadora de origen
86 df_train['Dataset'] = 'TRAIN'
87 df_test['Dataset'] = 'TEST'
88
89 # Unir datasets
90 df_combined = pd.concat([df_train, df_test], axis=0, ignore_index=True)
91
92 print(f"✓ Dataset combinado creado: {df_combined.shape}")
93 print(f" - Registros TRAIN: {(df_combined['Dataset'] == 'TRAIN').sum()}")
94 print(f" - Registros TEST: {(df_combined['Dataset'] == 'TEST').sum()}")
95
96 # Verificar integridad
97 print(f"\nVERIFICACIÓN INTEGRIDAD:")
98 print(f" - Nulos en Ahorro_final: {df_combined['Ahorro_final'].isnull().sum()}")
99 print(f" - Nulos en Mes_lici: {df_combined['Mes_lici'].isnull().sum()}")
100 print(f" - Rango Ahorro_final: [{df_combined['Ahorro_final'].min():.1f}%, {df_combined['Ahorro_final'].max():.1f}%]")
101
102 # =====
103 # PASO 4: PREPARACIÓN PARA CLUSTERING
104 # =====
105
106 print(f"\n4. PREPARACIÓN PARA CLUSTERING:")
107 print("-" * 50)
108
109 # Variables POST-licitación para clustering (5 variables especificadas)
110 variables_clustering = [
111     'Presupuesto_licitacion_lote_c', # Valor inicial de licitación
112     'N_ofertantes', # Número de ofertas recibidas
113     'Plazo_m', # Plazo de ejecución en meses
114     'Baja_p', # Porcentaje de baja económica
115     'C_precio_p' # Peso del criterio precio
116 ]
117
118 print(f"Variables para clustering:")
119 for i, var in enumerate(variables_clustering, 1):
120     print(f" {i}. {var}")
121
122 # Crear dataset para clustering
123 df_clustering = df_combined[variables_clustering + ['Ahorro_final', 'Dataset']].copy()
124 print(f"\nDataset clustering: {df_clustering.shape}")
125
126 # Verificar completitud y limpiar
127 nulos = df_clustering.isnull().sum()
128 print(f"Nulos por variable: {nulos.sum()} total")
129
130 df_clustering_clean = df_clustering.dropna()
131 print(f"Dataset limpio: {df_clustering_clean.shape}")
132
133 # =====
134 # PASO 5: NORMALIZACIÓN Y CLUSTERING
135 # =====
136
137 print(f"\n5. NORMALIZACIÓN Y CLUSTERING:")
138 print("-" * 50)

```

```

139
140 # Preparar datos para clustering (solo variables de entrada)
141 X_clustering = df_clustering_clean[variables_clustering].copy()
142 print(f"Matriz de clustering: {X_clustering.shape}")
143
144 # Normalizar datos
145 scaler = StandardScaler()
146 X_scaled = scaler.fit_transform(X_clustering)
147 print(f"✓ Datos normalizados con StandardScaler")
148
149 # Realizar clustering K=6 y K=8
150 kmeans_6 = KMeans(n_clusters=6, random_state=42, n_init=10)
151 kmeans_8 = KMeans(n_clusters=8, random_state=42, n_init=10)
152
153 clusters_6 = kmeans_6.fit_predict(X_scaled)
154 clusters_8 = kmeans_8.fit_predict(X_scaled)
155
156 # Agregar clusters al dataset limpio
157 df_clustering_clean['Cluster_6'] = clusters_6
158 df_clustering_clean['Cluster_8'] = clusters_8
159
160 print(f"✓ Clustering completado")
161 print(f"K=6 distribución: {np.bincount(clusters_6)}")
162 print(f"K=8 distribución: {np.bincount(clusters_8)}")
163
164 # Calcular métricas de calidad
165 silhouette_6 = silhouette_score(X_scaled, clusters_6)
166 silhouette_8 = silhouette_score(X_scaled, clusters_8)
167 inertia_6 = kmeans_6.inertia_
168 inertia_8 = kmeans_8.inertia_
169
170 print(f"\nMétricas de calidad:")
171 print(f"K=6 - Silhouette Score: {silhouette_6:.4f}, Inertia: {inertia_6:.2f}")
172 print(f"K=8 - Silhouette Score: {silhouette_8:.4f}, Inertia: {inertia_8:.2f}")
173
174 # =====
175 # PASO 6: INTEGRACIÓN AL DATASET COMPLETO
176 # =====
177
178 print(f"\n6. INTEGRACIÓN DE CLUSTERS AL DATASET COMPLETO:")
179 print("-" * 50)
180
181 # Crear dataset final con clusters
182 df_final = df_combined.copy()
183 df_final['Cluster_6'] = np.nan
184 df_final['Cluster_8'] = np.nan
185
186 # Asignar clusters usando índices
187 df_final.loc[df_clustering_clean.index, 'Cluster_6'] = df_clustering_clean['Cluster_6']
188 df_final.loc[df_clustering_clean.index, 'Cluster_8'] = df_clustering_clean['Cluster_8']
189
190 print(f"✓ Dataset final: {df_final.shape}")
191 print(f"✓ Cluster_6 asignados: {df_final['Cluster_6'].notna().sum()}")
192 print(f"✓ Cluster_8 asignados: {df_final['Cluster_8'].notna().sum()}")
193
194 # Guardar dataset final
195 df_final.to_excel('df_test_train_v8.xlsx', index=False)
196 print(f"✓ Dataset final guardado: df_test_train_v8.xlsx")
197
198 # =====
199 # PASO 7: ANÁLISIS ESTADÍSTICO POR CLUSTER
200 # =====
201
202 print(f"\n7. ANÁLISIS ESTADÍSTICO POR CLUSTER:")
203 print("-" * 50)
204
205 # Estadísticas K=6
206 cluster_6_stats = df_clustering_clean.groupby('Cluster_6')['Ahorro_final'].agg([
207     'count', 'mean', 'std', 'min', 'max'
208 ]).round(2)

```

```

209 cluster_6_stats.columns = ['Count', 'Mean_Ahorro', 'Std_Ahorro', 'Min_Ahorro',
210 'Max_Ahorro']
211 # Estadísticas K=8
212 cluster_8_stats = df_clustering_clean.groupby('Cluster_8')['Ahorro_final'].agg([
213     'count', 'mean', 'std', 'min', 'max'
214 ]).round(2)
215 cluster_8_stats.columns = ['Count', 'Mean_Ahorro', 'Std_Ahorro', 'Min_Ahorro',
216 'Max_Ahorro']
217 # Guardar estadísticas
218 cluster_6_stats.to_csv('estadisticas_k6_v8.csv')
219 cluster_8_stats.to_csv('estadisticas_k8_v8.csv')
220
221 print("✓ Estadísticas calculadas y guardadas")
222
223 # =====
224 # PASO 8: ANÁLISIS PCA PARA VISUALIZACIÓN
225 # =====
226
227 print(f"\n8. ANÁLISIS PCA PARA VISUALIZACIÓN:")
228 print("-" * 50)
229
230 # PCA para visualización
231 pca = PCA(n_components=2, random_state=42)
232 X_pca = pca.fit_transform(X_scaled)
233
234 print(f"✓ PCA completado")
235 print(f"Varianza explicada PC1: {pca.explained_variance_ratio_[0]:.4f}")
236 print(f"Varianza explicada PC2: {pca.explained_variance_ratio_[1]:.4f}")
237 print(f"Varianza total explicada: {pca.explained_variance_ratio_.sum():.4f}")
238
239 # Crear dataset para visualización
240 df_viz = pd.DataFrame({
241     'PC1': X_pca[:, 0],
242     'PC2': X_pca[:, 1],
243     'Cluster_6': clusters_6,
244     'Cluster_8': clusters_8,
245     'Ahorro_Final': df_clustering_clean['Ahorro_final'].values,
246     'Dataset': df_clustering_clean['Dataset'].values
247 })
248
249 df_viz.to_csv('datos_pca_v8.csv', index=False)
250 print("✓ Datos PCA guardados")
251
252 # =====
253 # PASO 9: ANÁLISIS DETALLADO POR CLUSTER
254 # =====
255
256 print(f"\n9. ANÁLISIS DETALLADO POR CLUSTER:")
257 print("-" * 50)
258
259 def analizar_cluster_detallado(df_data, cluster_col, cluster_id):
260     """Analiza un cluster específico con todas las variables"""
261
262     cluster_data = df_data[df_data[cluster_col] == cluster_id]
263
264     # Estadísticas básicas
265     count = len(cluster_data)
266     mean_ahorro = cluster_data['Ahorro_final'].mean()
267     std_ahorro = cluster_data['Ahorro_final'].std()
268     min_ahorro = cluster_data['Ahorro_final'].min()
269     max_ahorro = cluster_data['Ahorro_final'].max()
270
271     # Variables de clustering
272     mean_presupuesto = cluster_data['Presupuesto_licitacion_lote_c'].mean()
273     mean_ofertas = cluster_data['N_ofertantes'].mean()
274     mean_plazo = cluster_data['Plazo_m'].mean()
275     mean_baja = cluster_data['Baja_p'].mean()
276     mean_precio = cluster_data['C_precio_p'].mean()
277
278     # Distribución por dataset

```

```

279 train_pct = (cluster_data['Dataset'] == 'TRAIN').mean() * 100
280 test_pct = (cluster_data['Dataset'] == 'TEST').mean() * 100
281
282 # Categorización
283 if mean_ahorro < -7:
284     perfil = "Alto Sobrecoste"
285     riesgo = "MUY ALTO"
286 elif mean_ahorro < -3:
287     perfil = "Sobrecoste Moderado"
288     riesgo = "ALTO"
289 elif mean_ahorro < 2:
290     perfil = "Equilibrio"
291     riesgo = "MEDIO"
292 elif mean_ahorro < 8:
293     perfil = "Ahorro Moderado"
294     riesgo = "BAJO"
295 elif mean_ahorro < 12:
296     perfil = "Alto Ahorro"
297     riesgo = "MUY BAJO"
298 else:
299     perfil = "Ahorro Excepcional"
300     riesgo = "MUY BAJO"
301
302 return {
303     'Cluster': cluster_id,
304     'Perfil_Ahorro': perfil,
305     'Count': count,
306     'Pct_Total': count/len(df_data)*100,
307     'Mean_Ahorro': mean_ahorro,
308     'Std_Ahorro': std_ahorro,
309     'Mean_Presupuesto': mean_presupuesto,
310     'Mean_Oferatas': mean_ofertas,
311     'Mean_Plazo': mean_plazo,
312     'Mean_Baja': mean_baja,
313     'Mean_Precio': mean_precio,
314     'Train_Pct': train_pct,
315     'Test_Pct': test_pct,
316     'Riesgo': riesgo
317 }
318
319 # Análisis para K=6 y K=8
320 analisis_k6 = []
321 for i in range(6):
322     resultado = analizar_cluster_detallado(df_clustering_clean, 'Cluster_6', i)
323     analisis_k6.append(resultado)
324
325 analisis_k8 = []
326 for i in range(8):
327     resultado = analizar_cluster_detallado(df_clustering_clean, 'Cluster_8', i)
328     analisis_k8.append(resultado)
329
330 # Convertir a DataFrames y ordenar por ahorro
331 df_analisis_k6 = pd.DataFrame(analisis_k6).sort_values('Mean_Ahorro', ascending=False)
332 df_analisis_k8 = pd.DataFrame(analisis_k8).sort_values('Mean_Ahorro', ascending=False)
333
334 # Guardar análisis detallados
335 df_analisis_k6.to_csv('analisis_detallado_k6_v8.csv', index=False)
336 df_analisis_k8.to_csv('analisis_detallado_k8_v8.csv', index=False)
337
338 print(f"✓ Análisis detallado completado")
339
340 # =====
341 # PASO 10: RESUMEN FINAL
342 # =====
343
344 print(f"\n{'='*80}")
345 print("PROCESO COMPLETADO EXITOSAMENTE - V8")
346 print(f"{'='*80}")
347
348 print(f"\n■ RESUMEN EJECUTIVO:")
349 print(f"• Dataset combinado: {df_final.shape[0]} contratos ({df_final.shape[1]} variables)")

```

```
350 print(f"• Train: { (df_final['Dataset'] == 'TRAIN').sum() } contratos")
351 print(f"• Test: { (df_final['Dataset'] == 'TEST').sum() } contratos")
352 print(f"• Variables clustering: {len(variables_clustering)}")
353 print(f"• Ahorro promedio global: { df_final['Ahorro_final'].mean():.2f}%" )
354 print(f"• Sobrecostos totales: { (df_final['Ahorro_final'] < 0).mean()*100:.1f}%" )
355
356 print(f"\n🔗 MÉTRICAS CLUSTERING:")
357 print(f"K=6 - Silhouette: {silhouette_6:.4f}, Clusters exitosos: { (df_analisis_k6['Mean_Ahorro'] > 5).sum()}/6")
358 print(f"K=8 - Silhouette: {silhouette_8:.4f}, Clusters exitosos: { (df_analisis_k8['Mean_Ahorro'] > 5).sum()}/8")
359
360 print(f"\n✅ ARCHIVOS GENERADOS:")
361 print(f"1. df_test_train_v8.xlsx - Dataset combinado con clusters")
362 print(f"2. estadisticas_k6_v8.csv - Estadísticas clusters K=6")
363 print(f"3. estadisticas_k8_v8.csv - Estadísticas clusters K=8")
364 print(f"4. analisis_detallado_k6_v8.csv - Análisis completo K=6")
365 print(f"5. analisis_detallado_k8_v8.csv - Análisis completo K=8")
366 print(f"6. datos_pca_v8.csv - Datos para visualización PCA")
367
368 print(f"\n📌 LISTO PARA SIGUIENTE FASE:")
369 print(f"Random Forest con clusters como variable objetivo")
370 print(f"Variables PRE-licitación como predictores")
371
372 print(f"\n{'='*80}")
373 print("FIN DEL SCRIPT")
374 print(f"\n{'='*80}")
```