

Análisis de Datos con Python y Streamlit

Herramientas para Estadísticas Básicas, Limpieza de Datos, Análisis Visual Interactivo de Datos (ENAH) y Análisis de Gastos en Servicios Básicos con MANOVA

Edilberto Wilson Mamani Emanuel

4 de julio de 2025

1. Analizador de CSV - Estadísticas Básicas

Propósito: Esta herramienta permite cargar archivos CSV y realizar análisis estadísticos básicos de forma intuitiva. Está diseñada para usuarios que no tienen experiencia avanzada en estadística, proporcionando explicaciones claras de cada concepto.

Funcionalidades principales:

- Carga automática de archivos CSV con detección de codificación
- Cálculo de estadísticas descriptivas (media, mediana, moda, desviación estándar)
- Visualizaciones interactivas (histogramas, gráficos de caja)
- Interpretación automática de resultados
- Interfaz web amigable con Streamlit

Conceptos estadísticos incluidos:

- Medidas de tendencia central (promedio, mediana, moda)
- Medidas de dispersión (desviación estándar, rango, cuartiles)
- Análisis de distribución de datos
- Interpretación de coeficiente de variación

1.1. Código Fuente

```
1 import streamlit as st
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sn
6 from io import StringIO
7
8 # Configurar la página
9 st.set_page_config(
```

```

10     page_title="        Analizador de CSV - Estadísticas Básicas",
11     page_icon="        ",
12     layout="wide"
13 )
14
15 # Título principal
16 st.title("        Analizador de CSV - Estadísticas Básicas")
17 st.markdown("### Sube tu archivo CSV y descubre qué dicen tus datos!"
18 )
19
20 # Sidebar con información
21 st.sidebar.markdown("##        Guía Rápida")
22 st.sidebar.markdown("""
23 ** Qué puedes hacer aquí ?**
24 - Subir un archivo CSV
25 - Ver estadísticas básicas
26 - Entender qué significan los números
27 - Ver gráficos fáciles de entender
28 """)
29
30 # Función para explicar conceptos estadísticos
31 def explicar_concepto(concepto):
32     explicaciones = {
33         "promedio": "        **Promedio (Media)**: Es como dividir una
34         pizza entre todos. Si tienes 10 pizzas y 5 personas, cada uno come 2
35         pizzas en promedio.",
36         "mediana": "        **Mediana**: Es el valor del medio cuando
37         ordenas todos los números de menor a mayor. Como el estudiante que
38         queda en el medio en una fila ordenada por altura.",
39         "moda": "        **Moda**: Es el valor que más se repite. Como
40         el color favorito que eligi la mayoría de personas en una encuesta
41         .",
42         "desviacion": "        **Desviación Estándar**: Nos dice qué
43         tan dispersos están los datos. Si es pequeña, los datos están muy
44         juntos. Si es grande, están muy separados.",
45         "minimo": "        **Mínimo**: El valor más pequeño en tus
46         datos.",
47         "maximo": "        **Máximo**: El valor más grande en tus datos
48         .",
49         "rango": "        **Rango**: La diferencia entre el valor más
50         grande y el más pequeño.",
51         "cuartiles": "        **Cuartiles**: Dividen tus datos en 4
52         partes iguales. Como dividir a todos los estudiantes en 4 grupos del
53         mismo tamaño según sus calificaciones."
54     }
55     return explicaciones.get(concepto, "")
56
57 # Subir archivo
58 uploaded_file = st.file_uploader(
59     "        Sube tu archivo CSV aquí ",
60     type=['csv'],
61     help="Arrastra y suelta tu archivo CSV o haz clic para seleccionarlo"
62 )
63
64 if uploaded_file is not None:
65     try:
66         # Función para detectar la codificación del archivo

```

```

53     def leer_csv_con_codificacion(archivo):
54         # Lista de codificaciones comunes
55         codificaciones = ['utf-8', 'latin1', 'iso-8859-1', 'cp1252',
56                             'utf-16']
57
58         for codificacion in codificaciones:
59             try:
60                 # Resetear el puntero del archivo
61                 archivo.seek(0)
62                 # Intentar leer con la codificación actual
63                 df = pd.read_csv(archivo, encoding=codificacion)
64                 st.success(f"    Archivo leído correctamente con
65                             codificación: {codificacion}")
66                 return df
67             except UnicodeDecodeError:
68                 continue
69             except Exception as e:
70                 # Si hay otro error, intentar con la siguiente
71                 codificación
72                 continue
73
74         # Si ninguna codificación funciona, intentar con 'errors=
75         ignore'
76         try:
77             archivo.seek(0)
78             df = pd.read_csv(archivo, encoding='utf-8', errors='
79             ignore')
80             st.warning("    Se leyó el archivo ignorando algunos
81                             caracteres especiales")
82             return df
83         except:
84             raise Exception("No se pudo leer el archivo con ninguna
85                             codificación común")
86
87         # Leer el archivo CSV con detección automática de
88         codificación
89         df = leer_csv_con_codificacion(uploaded_file)
90
91         # Mostrar información básica del archivo
92         st.success(f"    Archivo cargado exitosamente!")
93
94         col1, col2, col3 = st.columns(3)
95         with col1:
96             st.metric("    Total de filas", len(df))
97         with col2:
98             st.metric("    Total de columnas", len(df.columns))
99         with col3:
100             st.metric("    Columnas numéricas", len(df.select_dtypes
101                             (include=[np.number]).columns))
102
103         # Mostrar vista previa de los datos
104         st.subheader("    Vista previa de tus datos")
105         st.dataframe(df.head(10), use_container_width=True)
106
107         # Obtener solo columnas numéricas
108         columnas_numericas = df.select_dtypes(include=[np.number]).
109         columns.tolist()

```

```

101         if len(columnas_numericas) > 0:
102             st.subheader("          Estadísticas Básicas")
103
104             # Selector de columna
105             columna_seleccionada = st.selectbox(
106                 "          Selecciona una columna para analizar:",
107                 columnas_numericas,
108                 help="Elige la columna de la que quieres ver las
estadísticas"
109             )
110
111             if columna_seleccionada:
112                 datos_columna = df[columna_seleccionada].dropna()
113
114                 # Calcular estadísticas
115                 promedio = datos_columna.mean()
116                 mediana = datos_columna.median()
117                 moda = datos_columna.mode().iloc[0] if not datos_columna
.mode().empty else "No hay moda"
118                 desviacion = datos_columna.std()
119                 minimo = datos_columna.min()
120                 maximo = datos_columna.max()
121                 rango = maximo - minimo
122                 q1 = datos_columna.quantile(0.25)
123                 q3 = datos_columna.quantile(0.75)
124
125                 # Mostrar estadísticas en tarjetas
126                 st.markdown(f"###          Estadísticas de la columna:
**{columna_seleccionada}**")
127
128                 # Fila 1 de métricas
129                 col1, col2, col3, col4 = st.columns(4)
130                 with col1:
131                     st.metric("          Promedio", f"{promedio:.2f}")
132                     st.info(explicar_concepto("promedio"))
133
134                 with col2:
135                     st.metric("          Mediana", f"{mediana:.2f}")
136                     st.info(explicar_concepto("mediana"))
137
138                 with col3:
139                     st.metric("          Mínimo", f"{minimo:.2f}")
140                     st.info(explicar_concepto("minimo"))
141
142                 with col4:
143                     st.metric("          Máximo", f"{maximo:.2f}")
144                     st.info(explicar_concepto("maximo"))
145
146                 # Fila 2 de métricas
147                 col1, col2, col3, col4 = st.columns(4)
148                 with col1:
149                     st.metric("          Desviación Estándar", f"{
desviacion:.2f}")
150                     st.info(explicar_concepto("desviacion"))
151
152                 with col2:
153                     st.metric("          Moda", f"{moda}")
154                     st.info(explicar_concepto("moda"))

```

```

155         with col3:
156             st.metric("          Rango", f"{rango:.2f}")
157             st.info(explicar_concepto("rango"))
158
159         with col4:
160             st.metric("          Datos v lidos", len(datos_columna
161 ))
162             st.info("Cantidad de datos que no est n vac os")
163
164         # Interpretaci n autom tica
165         st.subheader("          Qu   significa esto?")
166
167         interpretacion = []
168
169         if abs(promedio - mediana) / promedio < 0.1:
170             interpretacion.append("          **Datos balanceados**: Tu
promedio y mediana son muy similares, esto significa que tus datos
est n bien distribuidos.")
171         else:
172             interpretacion.append("          **Datos desbalanceados
**: Tu promedio y mediana son diferentes, esto puede indicar que
tienes valores muy altos o muy bajos que afectan el promedio.")
173
174         cv = (desviacion / promedio) * 100 if promedio != 0 else
0
175         if cv < 15:
176             interpretacion.append(f"          **Datos consistentes
**: Tus datos var an poco (variaci n del {cv:.1f}%), est n
bastante agrupados.")
177         elif cv < 30:
178             interpretacion.append(f"          **Datos moderadamente
variables**: Tus datos tienen variaci n media ({cv:.1f}%).")
179         else:
180             interpretacion.append(f"          **Datos muy variables
**: Tus datos est n muy dispersos (variaci n del {cv:.1f}%).")
181
182         for interp in interpretacion:
183             st.markdown(interp)
184
185         # Gr ficos
186         st.subheader("          Visualizaciones")
187
188         tab1, tab2, tab3 = st.tabs(["          Histograma", "
Gr fico de Caja", "          Estad sticas Resumidas"])
189
190         with tab1:
191             st.markdown("**   Qu   es un histograma?**          ")
192             st.markdown("Es como un gr fico de barras que
muestra cu ntos datos tienes en cada rango de valores. Te ayuda a
ver la forma de tus datos.")
193
194             fig, ax = plt.subplots(figsize=(10, 6))
195             ax.hist(datos_columna, bins=20, color='skyblue',
alpha=0.7, edgecolor='black')
196             ax.set_xlabel(columna_seleccionada)
197             ax.set_ylabel('Frecuencia (   Cuntas   veces aparece
?)')

```

```

198         ax.set_title(f'Distribuci n de {
columna_seleccionada}')
199         ax.grid(True, alpha=0.3)
200         st.pyplot(fig)
201
202     with tab2:
203         st.markdown("** Qu es un gr fico de caja?**
")
204         st.markdown("Es como un resumen visual de tus datos.
La caja muestra d nde est la mayor a de tus datos, y las l neas
muestran los valores extremos.")
205
206         fig, ax = plt.subplots(figsize=(10, 6))
207         ax.boxplot(datos_columna, vert=False)
208         ax.set_xlabel(columna_seleccionada)
209         ax.set_title(f'Resumen visual de {
columna_seleccionada}')
210         ax.grid(True, alpha=0.3)
211         st.pyplot(fig)
212
213     with tab3:
214         st.markdown("** Resumen completo de
estad sticas**")
215
216         # Crear tabla de resumen
217         resumen = pd.DataFrame({
218             'Estad stica': ['Promedio', 'Mediana', 'Moda',
'Desviaci n Est ndar', 'M nimo', 'M ximo', 'Rango', 'Cuartil 1
(25%)', 'Cuartil 3 (75%)'],
219             'Valor': [f"{promedio:.2f}", f"{mediana:.2f}",
str(moda), f"{desviacion:.2f}", f"{minimo:.2f}", f"{maximo:.2f}", f"{
rango:.2f}", f"{q1:.2f}", f"{q3:.2f}"],
220             'Significado': [
221                 'Valor t pico de tus datos',
222                 'Valor del medio cuando ordenas los datos',
223                 'Valor que m s se repite',
224                 'Qu tan dispersos est n los datos',
225                 'Valor m s peque o',
226                 'Valor m s grande',
227                 'Diferencia entre m ximo y m nimo',
228                 '25% de datos est n por debajo de este
valor',
229                 '75% de datos est n por debajo de este
valor'
230             ]
231         })
232
233         st.dataframe(resumen, use_container_width=True)
234
235         # An lisis de todas las columnas num ricas
236         if len(columnas_numericas) > 1:
237             st.subheader(" Resumen de todas las columnas
num ricas")
238
239             # Descripci n estad stica de todas las columnas
240             descripcion = df[columnas_numericas].describe()
241             st.dataframe(descripcion, use_container_width=True)
242

```

```

243         st.markdown("** Qu significa cada fila?**")
244         st.markdown("""
245         - **count**: Cu ntos datos v lidos tienes en cada
columna
246         - **mean**: El promedio de cada columna
247         - **std**: La desviaci n est ndar (dispersi n) de
cada columna
248         - **min**: El valor m nimo de cada columna
249         - **25%**: El 25% de los datos est n por debajo de
este valor
250         - **50%**: La mediana (valor del medio)
251         - **75%**: El 75% de los datos est n por debajo de
este valor
252         - **max**: El valor m ximo de cada columna
253         """)
254     else:
255         st.warning(" No se encontraron columnas num ricas en
tu archivo CSV. Las estad sticas solo se pueden calcular para datos
num ricos.")
256         st.markdown("** Qu puedes hacer?**")
257         st.markdown("- Verifica que tu archivo tenga columnas con
n meros ")
258         st.markdown("- Aseg rate de que los n meros no est n
escritos como texto")
259
260     except Exception as e:
261         st.error(f" Hubo un error al procesar tu archivo: {str(e)}
")
262         st.markdown("**Posibles soluciones:**")
263         st.markdown("- **Problema de codificaci n**: Abre tu archivo
CSV en un editor de texto y gu rdalo como UTF-8")
264         st.markdown("- **Archivo da ado**: Verifica que tu archivo CSV
no est corrupto")
265         st.markdown("- **Formato incorrecto**: Aseg rate de que sea un
archivo CSV v lido")
266         st.markdown("- **Caracteres especiales**: Evita usar caracteres
como , , , , , en los nombres de columnas")
267
268         # Mostrar informaci n adicional para debug
269         st.markdown("### Informaci n t cnica:")
270         st.code(f"Error espec fico: {str(e)}", language='text')
271
272         # Sugerencias espec ficas para codificaci n
273         st.markdown("### Si tienes problemas de codificaci n:")
274         st.markdown("""
275         **Opci n 1 - En Excel:**
276         1. Abre tu archivo en Excel
277         2. Ve a 'Archivo' 'Guardar como'
278         3. Selecciona 'CSV UTF-8 (delimitado por comas)'
279         4. Guarda el archivo
280
281         **Opci n 2 - En Google Sheets:**
282         1. Sube tu archivo a Google Sheets
283         2. Ve a 'Archivo' 'Descargar' 'Valores separados por
comas (.csv)'
284         3. Usa el archivo descargado
285         """)
286

```

```

287     st.markdown("###          Sigues teniendo problemas?")
288     st.markdown("Envíame las primeras líneas de tu archivo CSV y
    te ayudo a solucionarlo.")
289
290 else:
291     # Mostrar información de ayuda cuando no hay archivo
292     st.info("          Sube un archivo CSV para comenzar el análisis")
293
294     st.markdown("###          Qué son las estadísticas básicas?")
295     st.markdown("""
296     Las estadísticas básicas son como un resumen de tus datos que te
    ayuda a entender:
297
298     - **          Cul es el valor típico? (Promedio)
299     - **          Cul es el valor del medio? (Mediana)
300     - **          Cul es el valor más común? (Moda)
301     - **          Qué tan dispersos están los datos? (Desviación
    estándar)
302     - **          Cúales son los valores extremos? (Mínimo y
    máximo)
303     """)
304
305     st.markdown("###          Ejemplo de archivo CSV")
306     st.markdown("Tu archivo CSV debe verse algo así:")
307
308     ejemplo_csv = """Nombre,Edad,Salario,Experiencia
309 Juan,25,45000,2
310 María,30,55000,5
311 Pedro,35,65000,8
312 Ana,28,50000,3"""
313
314     st.code(ejemplo_csv, language='csv')
315
316     st.markdown("**          Consejos para tu archivo CSV:**")
317     st.markdown("- La primera fila debe contener los nombres de las
    columnas")
318     st.markdown("- Los datos deben estar separados por comas")
319     st.markdown("- Los números no deben contener símbolos como $ o %")
320     st.markdown("- Guarda el archivo con extensión .csv")
321
322 # Footer
323 st.markdown("----")
324 st.markdown("          **Consejo**: Las estadísticas te ayudan a tomar
    mejores decisiones basadas en tus datos. No necesitas ser un
    experto para entender lo que significan!")

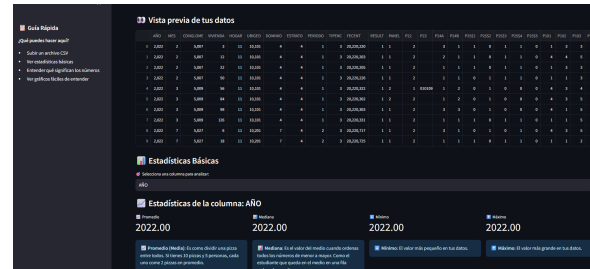
```

Listing 1: Analizador de CSV - Estadísticas Básicas

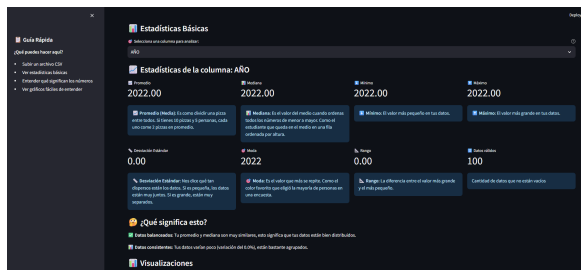
1.2. Capturas de Pantalla



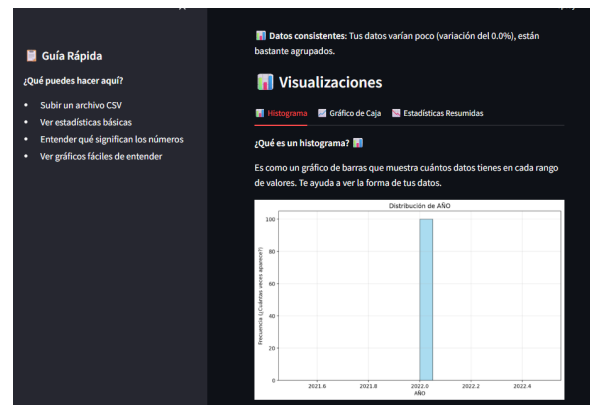
(a) Interfaz principal de carga de archivos



(b) Vista previa de datos cargados



(a) Métricas estadísticas calculadas



(b) Histograma de distribución

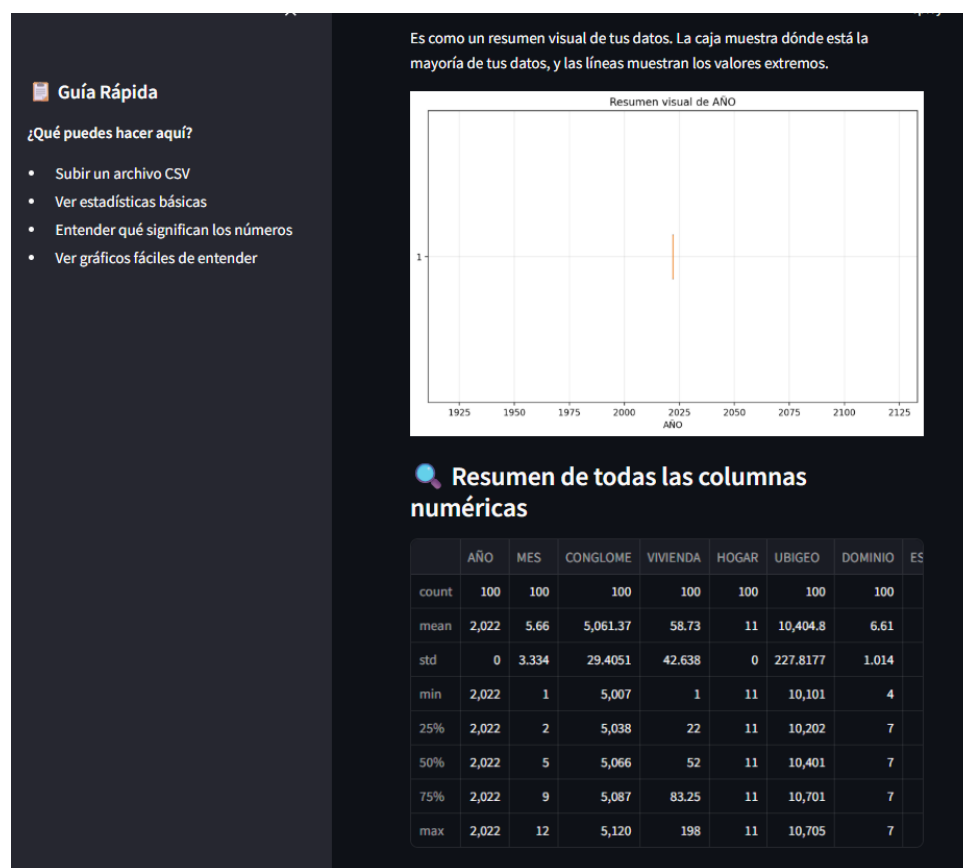


Figura 3: Gráfico de caja (boxplot) y resumen estadístico

1.3. Análisis Técnico

Aspectos técnicos destacados:

- **Detección automática de codificación:** El sistema prueba múltiples codificaciones (UTF-8, Latin-1, ISO-8859-1) para garantizar la correcta lectura de archivos CSV con caracteres especiales.
- **Interfaz adaptativa:** La aplicación se adapta automáticamente al tipo de datos encontrados, mostrando solo las opciones relevantes para cada conjunto de datos.
- **Interpretación automatizada:** El sistema compara automáticamente media y mediana para detectar asimetrías, y calcula el coeficiente de variación para evaluar la dispersión de los datos.
- **Visualizaciones interactivas:** Utiliza tanto Matplotlib como Plotly para crear gráficos informativos que ayudan a entender la distribución de los datos.

2. Limpieza de Datos y Detección de Outliers

Propósito: Esta herramienta se enfoca en la identificación y corrección de problemas comunes en conjuntos de datos, incluyendo valores faltantes, duplicados y valores atípicos (outliers). Es fundamental para preparar los datos antes del análisis.

Funcionalidades principales:

- Detección automática de datos faltantes con visualización de patrones
- Identificación de outliers usando métodos IQR y Z-Score
- Eliminación automática de duplicados
- Imputación de valores faltantes con múltiples estrategias
- Generación de reportes detallados de calidad de datos
- Descarga de datos limpios en formato CSV

Métodos de detección de outliers:

- **IQR (Rango Inter cuartílico):** Método robusto que usa Q1 y Q3 para definir límites
- **Z-Score:** Método basado en desviaciones estándar desde la media
- Visualizaciones con gráficos de caja y dispersión

2.1. Código Fuente

```

1 import streamlit as st
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from scipy import stats
7 import plotly.express as px
8 import plotly.graph_objects as go
9 from plotly.subplots import make_subplots
10 import warnings
11 warnings.filterwarnings('ignore')
12
13 # Configurar la página
14 st.set_page_config(
15     page_title="Limpieza de Datos y Detección de Outliers",
16     page_icon="📊",
17     layout="wide"
18 )
19
20 # Título principal
21 st.title("Limpieza de Datos y Detección de Outliers")
22 st.markdown("### Encuentra y limpia los datos problemáticos en tu archivo CSV!")
23
24 # Sidebar con información
25 st.sidebar.markdown("### Guía de Limpieza")
26 st.sidebar.markdown("""
27 ** ¿Qué hace esta herramienta? **
28 - Detecta datos faltantes
29 - Encuentra outliers (datos raros)
30 - Sugiere cómo limpiar tus datos
31 - Muestra gráficos fáciles de entender
32 - Genera reporte completo
33 """)

```

```

34
35 # Funci n para explicar conceptos
36 def explicar_concepto(concepto):
37     explicaciones = {
38         "outliers": "      **Outliers (Datos At picos)**: Son valores
que est n muy lejos del resto. Como una persona de 2.5 metros en un
grupo de personas normales.",
39         "datos_faltantes": "      **Datos Faltantes**": Son espacios
vac os en tu tabla. Como cuando no respondiste una pregunta en una
encuesta.",
40         "iqr": "      **M todo IQR**": Usa cuartiles para encontrar
datos raros. Si un dato est  muy lejos del 'grupo principal', lo
marca como extra o.",
41         "zscore": "      **M todo Z-Score**": Mide qu  tan 'normal' es
un dato. Si est  a m s de 3 desviaciones del promedio, es raro.",
42         "duplicados": "      **Datos Duplicados**": Filas que se repiten
exactamente igual. Como tener la misma persona dos veces en una
lista.",
43         "limpieza": "      **Limpieza de Datos**": Es como ordenar tu
cuarto. Quitas lo que no sirve y organizas lo que s  .
44     }
45     return explicaciones.get(concepto, "")
46
47 # Funci n para leer CSV con m ltiples codificaciones
48 def leer_csv_con_codificacion(archivo):
49     codificaciones = ['utf-8', 'latin1', 'iso-8859-1', 'cp1252', 'utf-16
',]
50
51     for codificacion in codificaciones:
52         try:
53             archivo.seek(0)
54             df = pd.read_csv(archivo, encoding=codificacion)
55             st.success(f"      Archivo le do correctamente con
codificaci n: {codificacion}")
56             return df
57         except UnicodeDecodeError:
58             continue
59         except Exception as e:
60             continue
61
62     try:
63         archivo.seek(0)
64         df = pd.read_csv(archivo, encoding='utf-8', errors='ignore')
65         st.warning("      Se ley  el archivo ignorando algunos
caracteres especiales")
66         return df
67     except:
68         raise Exception("No se pudo leer el archivo con ninguna
codificaci n com n")
69
70 # Funci n para detectar outliers con IQR
71 def detectar_outliers_iqr(df, columna):
72     Q1 = df[columna].quantile(0.25)
73     Q3 = df[columna].quantile(0.75)
74     IQR = Q3 - Q1
75     limite_inferior = Q1 - 1.5 * IQR
76     limite_superior = Q3 + 1.5 * IQR
77

```

```

78     outliers = df[(df[columna] < limite_inferior) | (df[columna] >
limite_superior)]
79     return outliers, limite_inferior, limite_superior
80
81 # Funci n para detectar outliers con Z-Score
82 def detectar_outliers_zscore(df, columna, umbral=3):
83     z_scores = np.abs(stats.zscore(df[columna].dropna()))
84     outliers = df[z_scores > umbral]
85     return outliers, umbral
86
87 # Funci n para generar reporte de limpieza
88 def generar_reporte_limpieza(df):
89     reporte = {}
90
91     # Informaci n general
92     reporte['filas_totales'] = len(df)
93     reporte['columnas_totales'] = len(df.columns)
94     reporte['celdas_totales'] = len(df) * len(df.columns)
95
96     # Datos faltantes
97     reporte['datos_faltantes'] = df.isnull().sum().sum()
98     reporte['porcentaje_faltantes'] = (reporte['datos_faltantes'] /
reporte['celdas_totales']) * 100
99
100    # Duplicados
101    reporte['filas_duplicadas'] = df.duplicated().sum()
102    reporte['porcentaje_duplicados'] = (reporte['filas_duplicadas'] /
reporte['filas_totales']) * 100
103
104    # Columnas con problemas
105    reporte['columnas_con_faltantes'] = df.isnull().sum()[df.isnull().
sum() > 0].to_dict()
106
107    return reporte
108
109 # Subir archivo
110 uploaded_file = st.file_uploader(
111     "        Sube tu archivo CSV aqu ",
112     type=['csv'],
113     help="Arrastra tu archivo CSV o haz clic para seleccionarlo"
114 )
115
116 if uploaded_file is not None:
117     try:
118         # Leer el archivo
119         df_original = leer_csv_con_codificacion(uploaded_file)
120         df = df_original.copy()
121
122         # Informaci n b sica
123         st.success("        Archivo        cargado exitosamente!")
124
125         # M tricas principales
126         col1, col2, col3, col4 = st.columns(4)
127         with col1:
128             st.metric("        Total de filas", len(df))
129         with col2:
130             st.metric("        Total de columnas", len(df.columns))
131         with col3:

```

```

132         st.metric("          Columnas num ricas", len(df.select_dtypes
(include=[np.number]).columns))
133         with col4:
134             st.metric("          Columnas de texto", len(df.select_dtypes(
include=['object']).columns))
135
136         # Generar reporte de limpieza
137         reporte = generar_reporte_limpieza(df)
138
139         # Mostrar vista previa
140         st.subheader("          Vista previa de tus datos")
141         st.dataframe(df.head(10), use_container_width=True)
142
143         # Tabs principales
144         tab1, tab2, tab3, tab4, tab5 = st.tabs([
145             "          Reporte General",
146             "          Datos Faltantes",
147             "          Outliers (Datos At picos)",
148             "          Limpieza Autom tica",
149             "          Reporte Final"
150         ])
151
152         with tab1:
153             st.subheader("          Reporte General de Calidad")
154
155             # M tricas de calidad
156             col1, col2 = st.columns(2)
157
158             with col1:
159                 st.metric(
160                     "          Datos Faltantes",
161                     f"{reporte['datos_faltantes']:,}",
162                     f"{reporte['porcentaje_faltantes']:.1f}% del total"
163                 )
164
165                 st.metric(
166                     "          Filas Duplicadas",
167                     f"{reporte['filas_duplicadas']:,}",
168                     f"{reporte['porcentaje_duplicados']:.1f}% del total"
169                 )
170
171             with col2:
172                 # Gr fico de calidad general
173                 fig = go.Figure(data=[
174                     go.Bar(
175                         name='Datos Completos',
176                         x=['Calidad de Datos'],
177                         y=[100 - reporte['porcentaje_faltantes']],
178                         marker_color='green'
179                     ),
180                     go.Bar(
181                         name='Datos Faltantes',
182                         x=['Calidad de Datos'],
183                         y=[reporte['porcentaje_faltantes']],
184                         marker_color='red'
185                     )
186                 ])
187

```

```

188         fig.update_layout(
189             title='Calidad General de los Datos',
190             yaxis_title='Porcentaje',
191             barmode='stack',
192             height=300
193         )
194         st.plotly_chart(fig, use_container_width=True)
195
196     # Calidad por columna
197     st.subheader("          Calidad por Columna")
198
199     calidad_columnas = []
200     for columna in df.columns:
201         faltantes = df[columna].isnull().sum()
202         porcentaje_faltantes = (faltantes / len(df)) * 100
203         porcentaje_completos = 100 - porcentaje_faltantes
204
205         calidad_columnas.append({
206             'Columna': columna,
207             'Datos Completos': int(len(df) - faltantes),
208             'Datos Faltantes': int(faltantes),
209             '% Completos': round(porcentaje_completos, 1),
210             '% Faltantes': round(porcentaje_faltantes, 1),
211             'Tipo de Dato': str(df[columna].dtype)
212         })
213
214     df_calidad = pd.DataFrame(calidad_columnas)
215     st.dataframe(df_calidad, use_container_width=True)
216
217     # Interpretaci n
218     st.info(explicar_concepto("datos_faltantes"))
219     st.info(explicar_concepto("duplicados"))
220
221     with tab2:
222         st.subheader("          An lisis de Datos Faltantes")
223
224         if reporte['datos_faltantes'] > 0:
225             # Mapa de calor de datos faltantes
226             st.markdown("###          Mapa de Calor de Datos Faltantes
227 ")
228             st.markdown("** Qu muestra? ** Las reas rojas son
229 donde faltan datos. Las azules donde est n completos.")
230
231             # Crear matriz de valores faltantes
232             missing_matrix = df.isnull().astype(int)
233
234             fig, ax = plt.subplots(figsize=(12, 8))
235             sns.heatmap(missing_matrix, cbar=True, cmap='RdYlBu_r',
236 ax=ax)
237             ax.set_title('Mapa de Datos Faltantes (Rojo = Falta,
238 Azul = Completo)')
239             st.pyplot(fig)
240
241             # Estad sticas por columna
242             st.markdown("###          Estad sticas Detalladas")
243
244             columnas_problematicas = []
245             for columna, faltantes in reporte['

```

```

columnas_con_faltantes'].items():
    porcentaje = (faltantes / len(df)) * 100

    # Determinar nivel de problema
    if porcentaje < 5:
        nivel = "M nimo"
    elif porcentaje < 20:
        nivel = "Moderado"
    elif porcentaje < 50:
        nivel = "Alto"
    else:
        nivel = "Cr tico"

    columnas_problematicas.append({
        'Columna': columna,
        'Datos Faltantes': faltantes,
        'Porcentaje': f"{porcentaje:.1f}%",
        'Nivel de Problema': nivel,
        'Recomendaci n': 'Eliminar columna' if
porcentaje > 70 else 'Rellenar valores' if porcentaje > 30 else '
Eliminar filas'
    })

df_problematicas = pd.DataFrame(columnas_problematicas)
st.dataframe(df_problematicas, use_container_width=True)

# Recomendaciones
st.markdown("### Recomendaciones para Datos
Faltantes")
st.markdown("""
- ** M nimo (< 5%):** Eliminar las filas con
datos faltantes
- ** Moderado (5-20%):** Rellenar con promedio,
mediana o moda
- ** Alto (20-50%):** Considerar t cnicas
avanzadas de imputaci n
- ** Cr tico (> 50%):** Evaluar si eliminar la
columna completa
""")

else:
    st.success("Excelente ! No tienes datos
faltantes en tu archivo.")

with tab3:
    st.subheader("Detecci n de Outliers (Datos At picos
)")

    st.info(explicar_concepto("outliers"))

    # Obtener columnas num ricas
    columnas_numericas = df.select_dtypes(include=[np.number]).
columns.tolist()

    if len(columnas_numericas) > 0:
        # Selector de columna
        columna_analizar = st.selectbox(
            "Selecciona una columna para analizar

```



```

outliers:",
289         columnas_numericas
290     )
291
292     if columna_analizar:
293         # Selector de m todo
294         metodo = st.radio(
295             "Selecciona el m todo de detecci n:",
296             ["IQR (Recomendado)", "Z-Score"],
297             help="IQR es m s robusto, Z-Score es m s
sensible"
298         )
299
300         col1, col2 = st.columns(2)
301
302         with col1:
303             if metodo == "IQR (Recomendado)":
304                 st.info(explicar_concepto("iqr"))
305
306                 # Detectar outliers con IQR
307                 outliers, limite_inf, limite_sup =
detectar_outliers_iqr(df, columna_analizar)
308
309                 st.metric("Outliers detectados", len(
outliers))
310                 st.metric("L mite inferior", f"{
limite_inf:.2f}")
311                 st.metric("L mite superior", f"{
limite_sup:.2f}")
312
313             else:
314                 st.info(explicar_concepto("zscore"))
315
316                 # Detectar outliers con Z-Score
317                 outliers, umbral = detectar_outliers_zscore(
df, columna_analizar)
318
319                 st.metric("Outliers detectados", len(
outliers))
320                 st.metric("Umbral Z-Score", f"{umbral
}")
321
322         with col2:
323             # Gr fico de caja
324             fig = px.box(df, y=columna_analizar, title=f'
Detecci n de Outliers: {columna_analizar}')
325             fig.update_layout(height=400)
326             st.plotly_chart(fig, use_container_width=True)
327
328             # Mostrar outliers detectados
329             if len(outliers) > 0:
330                 st.subheader("Outliers Detectados")
331
332             # Mostrar estad sticas de outliers
333             col1, col2, col3 = st.columns(3)
334             with col1:
335                 st.metric("Total de outliers", len(
outliers))

```

```

336         with col2:
337             st.metric("          Valor máximo outlier", f"
{outliers[columna_analizar].max():.2f}")
338         with col3:
339             st.metric("          Valor mínimo outlier", f"
{outliers[columna_analizar].min():.2f}")
340
341         # Tabla de outliers
342         st.dataframe(outliers, use_container_width=True)
343
344         # Gráfico de dispersión
345         st.subheader("          Visualización de Outliers"
)
346
347         # Crear gráfico de dispersión
348         fig = px.scatter(
349             df,
350             x=df.index,
351             y=columna_analizar,
352             title=f'Outliers en {columna_analizar}',
353             labels={'x': 'Índice de Fila', 'y':
columna_analizar}
354         )
355
356         # Destacar outliers
357         fig.add_scatter(
358             x=outliers.index,
359             y=outliers[columna_analizar],
360             mode='markers',
361             marker=dict(color='red', size=10),
362             name='Outliers'
363         )
364
365         st.plotly_chart(fig, use_container_width=True)
366
367         # Recomendaciones
368         st.subheader("          Qué hacer con los
outliers?")
369         porcentaje_outliers = (len(outliers) / len(df))
370         * 100
371
372         if porcentaje_outliers < 1:
373             st.success(f"          Solo {porcentaje_outliers
:.1f}% son outliers. Puedes eliminarlos de forma segura.")
374         elif porcentaje_outliers < 5:
375             st.warning(f"          {porcentaje_outliers:.1f
}% son outliers. Considera si son errores o datos válidos.")
376         else:
377             st.error(f"          {porcentaje_outliers:.1f}%
son outliers. Revisa si hay problemas en la recolección de datos.")
378
379         else:
380             st.success("          Excelente ! No se detectaron
outliers en esta columna.")
381
382         # Análisis de todas las columnas
383         st.subheader("          Resumen de Outliers en Todas las
Columnas")

```

```

383
384         resumen_outliers = []
385         for columna in columnas_numericas:
386             outliers_iqr, _ = detectar_outliers_iqr(df,
columna)
387             outliers_zscore, _ = detectar_outliers_zscore(df,
columna)
388
389             resumen_outliers.append({
390                 'Columna': columna,
391                 'Outliers IQR': len(outliers_iqr),
392                 'Outliers Z-Score': len(outliers_zscore),
393                 '% Outliers IQR': f"{(len(outliers_iqr)/len(df)
*100):.1f}%",
394                 '% Outliers Z-Score': f"{(len(outliers_zscore)/
len(df)*100):.1f}%"
395             })
396
397         df_resumen = pd.DataFrame(resumen_outliers)
398         st.dataframe(df_resumen, use_container_width=True)
399
400     else:
401         st.warning("          No se encontraron columnas num ricas
para analizar outliers.")
402
403     with tab4:
404         st.subheader("          Limpieza Autom tica de Datos")
405
406         st.info(explicar_concepto("limpieza"))
407
408         # Opciones de limpieza
409         st.markdown("###          Opciones de Limpieza")
410
411         col1, col2 = st.columns(2)
412
413         with col1:
414             st.markdown("**          Eliminar Datos**")
415             eliminar_duplicados = st.checkbox("Eliminar filas
duplicadas", value=True)
416             eliminar_filas_vacias = st.checkbox("Eliminar filas
completamente vac as", value=True)
417
418             # Umbral para eliminar filas con muchos faltantes
419             umbral_faltantes = st.slider(
420                 "Eliminar filas con m s de X% datos faltantes:",
421                 min_value=0, max_value=100, value=50, step=5
422             )
423
424         with col2:
425             st.markdown("**          Rellenar Datos Faltantes**")
426             rellenar_numericos = st.selectbox(
427                 "Rellenar columnas num ricas con:",
428                 ["No rellenar", "Promedio", "Mediana", "Moda", "Cero
"]
429             )
430
431             rellenar_texto = st.selectbox(
432                 "Rellenar columnas de texto con:",

```

```

433         ["No rellenar", "Moda", "Texto personalizado"]
434     )
435
436     if rellenar_texto == "Texto personalizado":
437         texto_personalizado = st.text_input("Texto para
rellenar:", value="Sin datos")
438
439     # Bot n para limpiar
440     if st.button("      Limpiar Datos", type="primary"):
441         df_limpio = df.copy()
442         acciones_realizadas = []
443
444         # Eliminar duplicados
445         if eliminar_duplicados:
446             duplicados_antes = df_limpio.duplicated().sum()
447             df_limpio = df_limpio.drop_duplicates()
448             if duplicados_antes > 0:
449                 acciones_realizadas.append(f"      Eliminadas {
duplicados_antes} filas duplicadas")
450
451         # Eliminar filas completamente vac as
452         if eliminar_filas_vacias:
453             filas_vacias = df_limpio.isnull().all(axis=1).sum()
454             df_limpio = df_limpio.dropna(how='all')
455             if filas_vacias > 0:
456                 acciones_realizadas.append(f"      Eliminadas {
filas_vacias} filas completamente vac as")
457
458         # Eliminar filas con muchos faltantes
459         if umbral_faltantes < 100:
460             umbral_decimal = umbral_faltantes / 100
461             filas_antes = len(df_limpio)
462             df_limpio = df_limpio.dropna(thresh=len(df_limpio.
columns) * (1 - umbral_decimal))
463             filas_eliminadas = filas_antes - len(df_limpio)
464             if filas_eliminadas > 0:
465                 acciones_realizadas.append(f"      Eliminadas {
filas_eliminadas} filas con m s del {umbral_faltantes}% de datos
faltantes")
466
467         # Rellenar datos num ricos
468         if rellenar_numericos != "No rellenar":
469             columnas_numericas = df_limpio.select_dtypes(include
=[np.number]).columns
470             for columna in columnas_numericas:
471                 faltantes_antes = df_limpio[columna].isnull().
sum()
472                 if faltantes_antes > 0:
473                     if rellenar_numericos == "Promedio":
474                         df_limpio[columna] = df_limpio[columna].
fillna(df_limpio[columna].mean())
475                     elif rellenar_numericos == "Mediana":
476                         df_limpio[columna] = df_limpio[columna].
fillna(df_limpio[columna].median())
477                     elif rellenar_numericos == "Moda":
478                         df_limpio[columna] = df_limpio[columna].
fillna(df_limpio[columna].mode().iloc[0] if not df_limpio[columna].
mode().empty else 0)

```

```

479         elif rellenar_numericos == "Cero":
480             df_limpio[columna] = df_limpio[columna].
fillna(0)
481
482             acciones_realizadas.append(f"           Rellenados
{faltantes_antes} valores faltantes en '{columna}' con {
rellenar_numericos.lower()}")
483
484             # Rellenar datos de texto
485             if rellenar_texto != "No rellenar":
486                 columnas_texto = df_limpio.select_dtypes(include=[
object']).columns
487                 for columna in columnas_texto:
488                     faltantes_antes = df_limpio[columna].isnull().
sum()
489                     if faltantes_antes > 0:
490                         if rellenar_texto == "Moda":
491                             df_limpio[columna] = df_limpio[columna].
fillna(df_limpio[columna].mode().iloc[0] if not df_limpio[columna].
mode().empty else "Sin datos")
492                         elif rellenar_texto == "Texto personalizado"
:
493                             df_limpio[columna] = df_limpio[columna].
fillna(texto_personalizado)
494
495             acciones_realizadas.append(f"           Rellenados
{faltantes_antes} valores faltantes en '{columna}')"
496
497             # Mostrar resultados
498             st.subheader("           Resultados de la Limpieza")
499
500             if acciones_realizadas:
501                 for accion in acciones_realizadas:
502                     st.write(accion)
503
504             # Comparaci n antes/despu s
505             col1, col2 = st.columns(2)
506
507             with col1:
508                 st.markdown("**           Antes de la limpieza:**")
509                 st.metric("Filas", len(df))
510                 st.metric("Datos faltantes", df.isnull().sum().
sum())
511                 st.metric("Duplicados", df.duplicated().sum())
512
513             with col2:
514                 st.markdown("**           Despu s de la limpieza:**")
515                 st.metric("Filas", len(df_limpio))
516                 st.metric("Datos faltantes", df_limpio.isnull().
sum().sum())
517                 st.metric("Duplicados", df_limpio.duplicated().
sum())
518
519             # Mostrar datos limpios
520             st.subheader("           Datos Limpios")
521             st.dataframe(df_limpio.head(10), use_container_width
=True)
522

```

```

523         # Guardar datos limpios en session state
524         st.session_state['datos_limpios'] = df_limpio
525
526         # Bot n para descargar
527         csv_limpio = df_limpio.to_csv(index=False)
528         st.download_button(
529             label="Descargar Datos Limpios",
530             data=csv_limpio,
531             file_name="datos_limpios.csv",
532             mime="text/csv"
533         )
534
535         else:
536             st.info("No se realizaron cambios. Tus datos
ya est n limpios.")
537
538     with tab5:
539         st.subheader("Reporte Final de Limpieza")
540
541         # Generar reporte final
542         st.markdown("### Resumen Ejecutivo")
543
544         # Crear m tricas de resumen
545         col1, col2, col3, col4 = st.columns(4)
546
547         with col1:
548             calidad_general = 100 - reporte['porcentaje_faltantes']
549             st.metric("Calidad General", f"{calidad_general
:.1f}%")
550
551         with col2:
552             problemas_totales = reporte['datos_faltantes'] + reporte
['filas_duplicadas']
553             st.metric("Problemas Detectados", f"{
problemas_totales:,}")
554
555         with col3:
556             columnas_problematicas = len(reporte['
columnas_con_faltantes'])
557             st.metric("Columnas con Problemas",
columnas_problematicas)
558
559         with col4:
560             if 'datos_limpios' in st.session_state:
561                 mejora = len(df) - len(st.session_state['
datos_limpios'])
562                 st.metric("Filas Eliminadas", mejora)
563             else:
564                 st.metric("Filas Eliminadas", "No aplicado")
565
566         # Reporte detallado
567         st.markdown("### Reporte Detallado")
568
569         reporte_detallado = f"""
570         ** REPORTE DE CALIDAD DE DATOS**
571
572         ** Informaci n General:**
573         - Total de filas: {reporte['filas_totales']:,}

```

```

574         - Total de columnas: {reporte['columnas_totales']:,}
575         - Total de celdas: {reporte['celdas_totales']:,}
576
577     **      Problemas Detectados:**
578     - Datos faltantes: {reporte['datos_faltantes']:,} ({reporte
['porcentaje_faltantes']:.1f}%)
579     - Filas duplicadas: {reporte['filas_duplicadas']:,} ({
reporte['porcentaje_duplicados']:.1f}%)
580
581     **      An lisis por Columna:**
582     """
583
584     for columna, faltantes in reporte['columnas_con_faltantes'].
items():
585         porcentaje = (faltantes / len(df)) * 100
586         reporte_detallado += f"\n- {columna}: {faltantes}
faltantes ({porcentaje:.1f}%"
587
588         if len(reporte['columnas_con_faltantes']) == 0:
589             reporte_detallado += "\n-          No hay columnas con datos
faltantes"
590
591     # Recomendaciones
592     reporte_detallado += f"""
593
594     **      Recomendaciones:**
595     """
596
597     if reporte['porcentaje_faltantes'] < 5:
598         reporte_detallado += "\n-          Calidad excelente.
M nimos problemas detectados."
599     elif reporte['porcentaje_faltantes'] < 20:
600         reporte_detallado += "\n-          Calidad buena.
Considerar limpiar datos faltantes."
601     elif reporte['porcentaje_faltantes'] < 50:
602         reporte_detallado += "\n-          Calidad regular.
Limpieza necesaria antes del an lisis."
603     else:
604         reporte_detallado += "\n-          Calidad cr tica.
Revisar proceso de recolecci n de datos."
605
606     if reporte['filas_duplicadas'] > 0:
607         reporte_detallado += f"\n-          Eliminar {reporte['
filas_duplicadas']} filas duplicadas"
608
609     if len(columnas_numericas) > 0:
610         reporte_detallado += "\n-          Revisar outliers en
columnas num ricas"
611
612     st.markdown(reporte_detallado)
613
614     # Gr fico de resumen
615     st.markdown("###          Visualizaci n del Reporte")
616
617     # Crear gr fico de pastel para mostrar distribuci n de
problemas
618     fig = go.Figure(data=[go.Pie(
619         labels=['Datos Completos', 'Datos Faltantes', '

```

```

Duplicados'],
620         values=[
621             reporte['celdas_totales'] - reporte['datos_faltantes
'] - reporte['filas_duplicadas'],
622             reporte['datos_faltantes'],
623             reporte['filas_duplicadas']
624         ],
625         hole=0.3,
626         marker_colors=['green', 'red', 'orange']
627     )))
628
629     fig.update_layout(
630         title="Distribuci n de Calidad de Datos",
631         annotations=[dict(text='Calidad<br>General', x=0.5, y
=0.5, font_size=16, showarrow=False)]
632     )
633
634     st.plotly_chart(fig, use_container_width=True)
635
636     # Descargar reporte
637     st.download_button(
638         label="Descargar Reporte Completo",
639         data=reporte_detallado,
640         file_name="reporte_limpieza_datos.txt",
641         mime="text/plain"
642     )
643
644     # Conclusiones finales
645     st.markdown("### Conclusiones y Pr ximos Pasos")
646
647     conclusiones = []
648
649     if reporte['porcentaje_faltantes'] < 5 and reporte['
filas_duplicadas'] == 0:
650         conclusiones.append(" ** Excelente !** Tus datos
est n en muy buena condici n.")
651     elif reporte['porcentaje_faltantes'] < 20:
652         conclusiones.append(" **Buena calidad** general, solo
necesitas limpieza menor.")
653     else:
654         conclusiones.append(" **Necesitas limpieza** antes
de hacer an lisis importantes.")
655
656     if len(columnas_numericas) > 0:
657         conclusiones.append(" **Revisa los outliers** en
la pesta a correspondiente.")
658
659     if reporte['filas_duplicadas'] > 0:
660         conclusiones.append(" **Elimina los duplicados
** para evitar sesgos en tu an lisis.")
661
662     conclusiones.append(" **Usa la limpieza autom tica**
para mejorar la calidad de tus datos.")
663
664     for conclusion in conclusiones:
665         st.markdown(conclusion)
666
667     # Recordatorio importante

```



```

668         st.warning("**Recordatorio importante**: Siempre
revisa manualmente los cambios propuestos antes de aplicarlos a tus
datos originales.")
669
670     except Exception as e:
671         st.error(f"Hubo un error al procesar tu archivo: {str(e)}")
672
673         st.markdown("**Posibles soluciones:**")
674         st.markdown("- **Problema de codificaci n**: Guarda tu archivo
CSV como UTF-8")
675         st.markdown("- **Archivo da ado**: Verifica que tu archivo CSV
no est corrupto")
676         st.markdown("- **Formato incorrecto**: Aseg rate de que sea un
archivo CSV v lido")
677         st.markdown("- **Caracteres especiales**: Evita caracteres
especiales en nombres de columnas")
678
679         # Informaci n t cnica
680         st.markdown("### Informaci n t cnica:")
681         st.code(f"Error espec fico: {str(e)}", language='text')
682     else:
683         # P gina de inicio cuando no hay archivo
684         st.info("Sube un archivo CSV para comenzar el an lisis de
limpieza")
685
686         st.markdown("### Qu son los datos sucios?")
687         st.markdown("""
688 Los datos "sucios" son datos que tienen problemas como:
689
690 - ** Datos faltantes**: Celdas vac as o con valores nulos
691 - ** Duplicados**: Filas que se repiten exactamente
692 - ** Outliers**: Valores muy diferentes al resto (datos
at picos)
693 - ** Inconsistencias**: Diferentes formatos para el mismo tipo
de dato
694 - ** Errores de tipeo**: Nombres mal escritos o caracteres
raros
695 """)
696
697         st.markdown("### Por qu limpiar los datos?")
698         st.markdown("""
699 - ** An lisis m s precisos**: Resultados m s confiables
700 - ** Mejores decisiones**: Conclusiones basadas en datos
correctos
701 - ** Procesamiento m s r pido**: Menos datos problem ticos
702 - ** Patrones m s claros**: Tendencias m s f ciles de
identificar
703 """)
704
705         st.markdown("### Ejemplo de datos que necesitan limpieza:")
706
707         ejemplo_sucio = """Nombre,Edad,Salario,Ciudad
708 Juan,25,45000,Madrid
709 Mar a ,,55000,Barcelona
710 Pedro,35,65000,Madrid
711 Ana,28,50000,
712 Juan,25,45000,Madrid

```

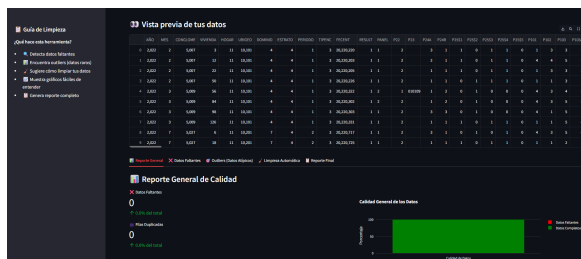
```

713 Carlos,200,70000,Valencia"""
714
715 st.code(ejemplo_sucio, language='csv')
716
717 st.markdown("**          Problemas en este ejemplo:**")
718 st.markdown("- Edad faltante para Mar a")
719 st.markdown("- Ciudad faltante para Ana")
720 st.markdown("- Juan est duplicado")
721 st.markdown("- Carlos tiene 200 a os (posible outlier)")
722
723 st.markdown("###          Lo que esta herramienta hace por ti:")
724 st.markdown("""
725 - **          Detecta autom ticamente** todos los problemas
726 - **          Muestra gr ficos** f ciles de entender
727 - **          Limpia autom ticamente** con un clic
728 - **          Genera reportes** detallados
729 - **          Da recomendaciones** espec ficas
730 """)
731
732 # Footer
733 st.markdown("----")
734 st.markdown("**Recuerda**: Datos limpios = An lisis confiables =
Mejores decisiones")

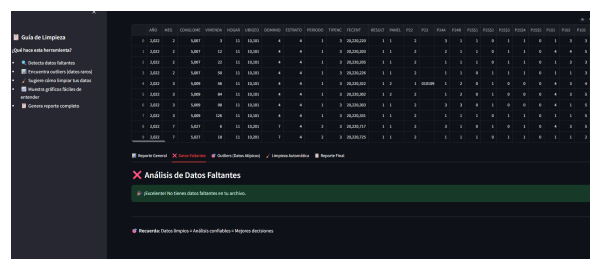
```

Listing 2: Limpieza de Datos y Detección de Outliers

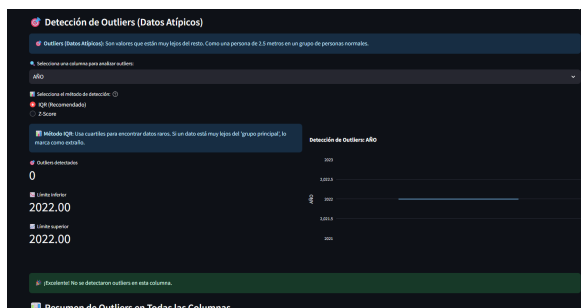
2.2. Capturas de Pantalla



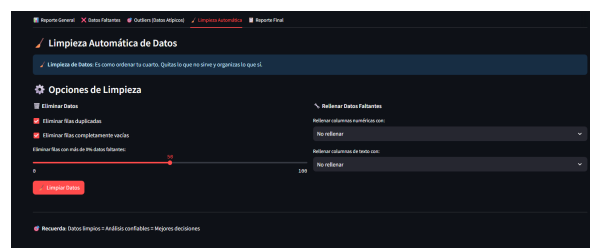
(a) Reporte general de calidad de datos



(b) Mapa de calor de datos faltantes



(a) Detección de outliers con método IQR



(b) Opciones de limpieza automática

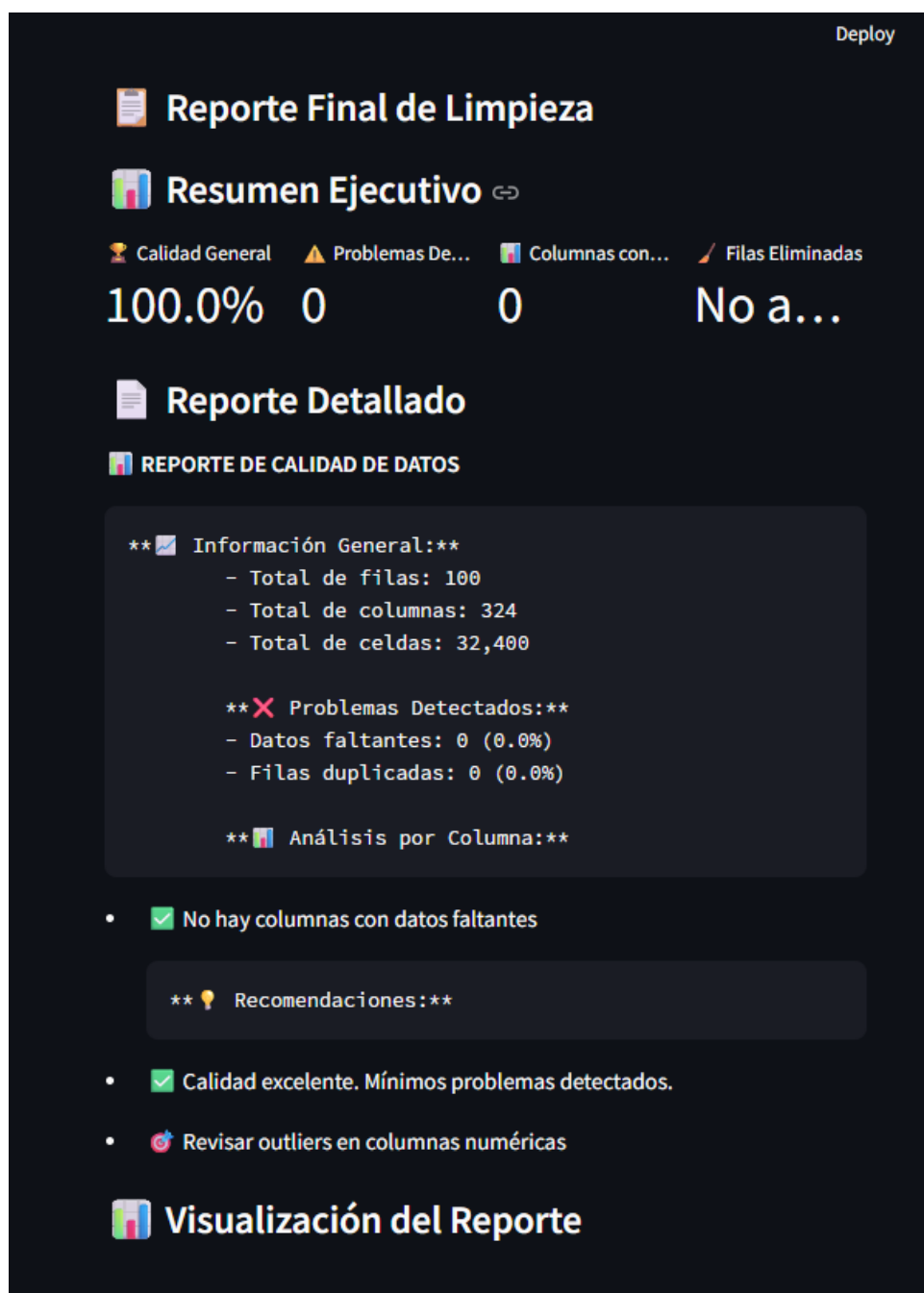


Figura 6: Reporte final con métricas de mejora

2.3. Análisis Técnico

Aspectos técnicos destacados:

- **Detección de outliers múltiple:** Implementa tanto el método IQR como Z-Score, permitiendo al usuario elegir el más apropiado según sus necesidades.
- **Visualización avanzada:** Utiliza Plotly para crear gráficos interactivos que permiten explorar los datos de forma dinámica, incluyendo mapas de calor para patrones de datos faltantes.
- **Estrategias de imputación:** Ofrece múltiples opciones para el tratamiento de

valores faltantes: eliminación, imputación por media/mediana/moda, o valores personalizados.

- **Análisis de calidad:** Calcula métricas de calidad de datos y genera reportes automáticos con recomendaciones específicas basadas en los problemas encontrados.
- **Preservación de datos:** Mantiene los datos originales intactos y trabaja con copias, permitiendo comparaciones antes/después de la limpieza.

3. Análisis Visual Interactivo de Datos (ENAH0)

Propósito: Facilitar la exploración y comprensión de conjuntos de datos CSV (especialmente los de encuestas como ENAH0) mediante visualizaciones interactivas, permitiendo a usuarios de cualquier nivel técnico realizar análisis exploratorios sin necesidad de programar.

Funcionalidades principales:

- **Carga robusta de archivos CSV:** Sube y lee archivos CSV automáticamente, gestionando diversas codificaciones de caracteres para evitar errores.
- **Normalización de nombres de columnas:** Limpia y estandariza los nombres de las columnas para un manejo más sencillo y consistente.
- **Matriz de correlación dinámica:** Genera un mapa de calor interactivo que muestra las relaciones lineales entre las variables numéricas seleccionadas.
- **Análisis de tendencias y distribuciones por categoría:** Visualiza cómo una variable numérica se comporta a través de diferentes grupos (ej. estrato social o región) usando gráficos de línea y de caja.
- **Mapa de calor de promedios combinados:** Explora el promedio de una variable numérica para cada combinación de dos categorías distintas, identificando patrones clave.

Características técnicas:

- **Framework Streamlit:** Desarrollado sobre Streamlit para una interfaz de usuario web intuitiva y fácil de usar.
- **Manipulación de datos con Pandas:** Utiliza la potente librería Pandas para el procesamiento, limpieza y análisis de los DataFrames.
- **Visualizaciones interactivas con Plotly Express:** Emplea Plotly para crear gráficos dinámicos y de alta calidad que facilitan la exploración de los datos.
- **Manejo integral de errores:** Incorpora validaciones y mensajes claros para guiar al usuario ante problemas de carga o selección de datos.

3.1. Código Fuente

```

1 import streamlit as st
2 import pandas as pd
3 import plotly.express as px
4 import plotly.graph_objects as go
5 import numpy as np
6 import io
7
8 st.set_page_config(page_title="Análisis ENAHO 2022", layout="wide")
9 st.title("Análisis Visual de Datos ENAHO")
10
11 # Subida de archivo
12 archivo = st.file_uploader("Sube tu archivo CSV", type=["csv"])
13
14 df = None # Inicializar df fuera del bloque if
15
16 if archivo is not None:
17     try:
18         # Intentar cargar con codificaciones comunes
19         codificaciones = ['utf-8', 'latin1', 'iso-8859-1', 'cp1252']
20         codificacion_exitosa = None
21         for encoding in codificaciones:
22             try:
23                 archivo.seek(0)
24                 df = pd.read_csv(archivo, encoding=encoding, low_memory=
False)
25                 codificacion_exitosa = encoding
26                 st.success(f"Datos cargados correctamente con
codificación: {encoding}")
27                 break
28             except UnicodeDecodeError:
29                 continue
30             except Exception as e:
31                 st.error(f"Error al leer con {encoding}: {e}")
32                 df = None
33
34         if df is None:
35             st.error("No se pudo cargar el archivo CSV con ninguna
de las codificaciones intentadas. Por favor, verifica el formato y la
codificación de tu archivo.")
36             st.stop()
37
38         st.dataframe(df.head())
39
40         # --- Limpiar los nombres de las columnas para mayor robustez
41         ---
42         original_columns = df.columns.tolist()
43         df.columns = df.columns.str.strip().str.replace(' ', '_').str.
replace('.', '_').str.replace('-', '_')
44
45         if isinstance(df.columns, pd.MultiIndex):
46             df.columns = df.columns.map('_'.join)
47             st.info("Detectado y simplificado un índice de columnas
multi-nivel.")
48
49         if original_columns != df.columns.tolist():
50             st.info("Se han limpiado los nombres de las columnas. Los

```

```

nuevos nombres son: " + ", ".join(df.columns.tolist())

# Selecci n de todas las columnas num ricas disponibles
todas_columnas_numericas = df.select_dtypes(include=np.number).
columns.tolist()

# --- 1. Matriz de correlaci n ---
st.subheader("1          Matriz de correlaci n de Variables
Num ricas Relevantes")
st.markdown("Selecciona las variables num ricas que consideres
m s relevantes para analizar su correlaci n.")

# Nuevo: Selector m ltiple para elegir las columnas a
correlacionar
selected_numeric_columns = st.multiselect(
    "Elige las columnas num ricas para la matriz de
correlaci n:",
    todas_columnas_numericas,
    default=todas_columnas_numericas[:5] if len(
todas_columnas_numericas) > 0 else [] # Selecciona las primeras 5 por
defecto si existen
)

if len(selected_numeric_columns) >= 2:
    corr = df[selected_numeric_columns].corr()

    fig_corr = px.imshow(
        corr,
        text_auto=".2f", # Formato a 2 decimales
        aspect="auto",
        color_continuous_scale='RdBu_r', # Escala de color Rojo-
Azul invertida
        title="Mapa de Correlaci n entre Variables Num ricas
Seleccionadas"
    )

    # Ajustar el tama o del texto dentro de las celdas y las
etiquetas de los ejes
    # El tama o de la fuente se har m s grande si hay menos
columnas
    font_size_labels = 12 if len(selected_numeric_columns) < 15
else 8
    font_size_text = 12 if len(selected_numeric_columns) < 10
else 8

    fig_corr.update_traces(textfont_size=font_size_text)
    fig_corr.update_xaxes(tickangle=-45, tickfont_size=
font_size_labels)
    fig_corr.update_yaxes(tickfont_size=font_size_labels)

    # Ajustar layout para mejor visualizaci n
    fig_corr.update_layout(
        height=max(500, len(selected_numeric_columns) * 40), #
Altura din mica
        width=max(500, len(selected_numeric_columns) * 40), #
Ancho din mico
        title_x=0.5,
        title_font_size=20

```

```

91         )
92
93         st.plotly_chart(fig_corr, use_container_width=True)
94
95         st.markdown("""
96         **Interpretaci n de la Matriz de Correlaci n:**
97         - **Valores cercanos a 1:** Indican una fuerte correlaci n
positiva (cuando una variable aumenta, la otra tambi n tiende a
aumentar).
98         - **Valores cercanos a -1:** Indican una fuerte correlaci n
negativa (cuando una variable aumenta, la otra tiende a disminuir).
99         - **Valores cercanos a 0:** Indican poca o ninguna
correlaci n lineal entre las variables.
100         """)
101
102     else:
103         st.warning("Por favor, selecciona al menos 2 columnas
num ricas para generar la matriz de correlaci n.")
104
105         # --- Validaci n y selecci n de la columna de resultado para
los siguientes gr ficos ---
106         available_columns = df.columns.tolist()
107
108         # Filtrar solo columnas num ricas para el resultado por defecto
, si es posible
109         opciones_resultado = [col for col in available_columns if pd.api
.types.is_numeric_dtype(df[col])]
110
111         if not opciones_resultado:
112             st.error("No se encontraron columnas num ricas para usar
como 'Resultado de la entrevista'. Algunos gr ficos no podr n
generarse.")
113             col_resultado_comun = None
114         else:
115             col_resultado_comun = st.selectbox(
116                 "Selecciona la columna de **Resultado de la entrevista**
(para gr ficos 2, 3 y 4)",
117                 opciones_resultado, # Limitar a columnas num ricas
key="resultado_select_comun"
118             )
119
120             # Asegurarse de que la columna seleccionada sigue siendo
num rica (por si el usuario elige otra cosa en un set distinto)
121             if not pd.api.types.is_numeric_dtype(df[col_resultado_comun
]):
122                 st.warning("La columna de Resultado seleccionada no es
num rica. Por favor, selecciona una columna num rica.")
123                 col_resultado_comun = None # Invalidar si no es
num rica
124
125             if col_resultado_comun: # Solo proceder si hay una columna de
resultado num rica v lida
126                 # --- 2. Gr fico de l nea: Estrato vs Resultado ---
127                 st.subheader("2          Influye el estrato social en el
resultado?")
128
129                 col_estrato = st.selectbox("Selecciona la columna de Estrato
Social", available_columns, key="estrato_select")
130

```

```

131         st.info(f"DEBUG: Columna 'Estrato Social' seleccionada: '{col_estrato}'")
132         st.info(f"DEBUG: Tipo de dato de '{col_estrato}': {df[col_estrato].dtype}")
133         st.info(f"DEBUG: Columna 'Resultado' seleccionada: '{col_resultado_comun}'")
134         st.info(f"DEBUG: Tipo de dato de '{col_resultado_comun}': {df[col_resultado_comun].dtype}")
135
136         if col_estrato not in df.columns:
137             st.error("La columna de estrato seleccionada no es v lida.")
138         else:
139             if not pd.api.types.is_numeric_dtype(df[col_estrato])
140 and not pd.api.types.is_categorical_dtype(df[col_estrato]):
141                 try:
142                     df[col_estrato] = df[col_estrato].astype('category')
143                     st.info(f"DEBUG: Columna '{col_estrato}' convertida a tipo 'category' para agrupaci n.")
144                 except Exception as e:
145                     st.error(f"Error al intentar convertir '{col_estrato}' a categor a: {e}. Esto podr a causar problemas en la agrupaci n.")
146
147             df_line = df[[col_estrato, col_resultado_comun]].dropna()
148
149             if df_line.empty:
150                 st.warning(f"No hay datos disponibles para '{col_estrato}' y '{col_resultado_comun}' despu s de eliminar valores faltantes.")
151             else:
152                 if df_line[col_estrato].nunique() < 2:
153                     st.warning(f"La columna '{col_estrato}' tiene muy pocos valores nicos ({df_line[col_estrato].nunique()}) para un an lisis significativo. Necesitas al menos 2 grupos.")
154                 else:
155                     try:
156                         st.info(f"DEBUG: Tipo de df_line['{col_estrato}'] antes de groupby: {df_line[col_estrato].dtype}")
157                         st.info(f"DEBUG: Shape de df_line['{col_estrato}'] antes de groupby: {df_line[col_estrato].shape}")
158
159                         df_grouped = df_line.groupby(col_estrato)[col_resultado_comun].mean().reset_index()
160
161                         fig2 = px.line(df_grouped, x=col_estrato, y=col_resultado_comun, markers=True, title=f"Tendencia del promedio de {col_resultado_comun} seg n {col_estrato}")
162
163                         st.plotly_chart(fig2, use_container_width=True)
164
165                         st.markdown(f"""
166                             **Interpretaci n:** Este gr fico de l nea muestra c mo el valor promedio de **'{col_resultado_comun}'** var a a trav s de las diferentes

```



```

categorías de **'{col_estrato}'**.
166         Puedes observar tendencias o patrones en los
        resultados a medida que cambias de estrato.
167     """
168     except Exception as e:
169         st.error(f"Error al generar el gráfico de
        línea: {e}. Esto podría deberse a datos inesperados en la columna
        de agrupación, o si la columna '{col_estrato}' tiene una estructura
        compleja (ej. MultiIndex).")
170         st.info(f"Asegúrate de que la columna '{
        col_estrato}' sea un tipo de dato simple y 1-dimensional (categórico
        , texto, o numérico discreto).")
171
172
173     # --- 3. Gráfico de caja: Región vs Resultado ---
174     st.subheader("3 Hay diferencias por región?")
175     col_region = st.selectbox("Selecciona la columna de Región"
        , available_columns, key="region_select")
176
177     if col_region not in df.columns:
178         st.error("La columna de región seleccionada no es
        válida.")
179     else:
180         if not pd.api.types.is_numeric_dtype(df[col_region]) and
        not pd.api.types.is_categorical_dtype(df[col_region]):
181             try:
182                 df[col_region] = df[col_region].astype('category
        ')
183                 st.info(f"DEBUG: Columna '{col_region}'
        convertida a tipo 'category' para agrupación.")
184             except Exception as e:
185                 st.error(f"Error al intentar convertir '{
        col_region}' a categoría: {e}. Esto podría causar problemas en la
        agrupación.")
186
187         df_box = df[[col_region, col_resultado_comun]].dropna()
188
189         if df_box.empty:
190             st.warning(f"No hay datos disponibles para '{
        col_region}' y '{col_resultado_comun}' después de eliminar valores
        faltantes.")
191         else:
192             if df_box[col_region].nunique() < 2:
193                 st.warning(f"La columna '{col_region}' tiene muy
        pocos valores únicos ({df_box[col_region].nunique()}) para un
        análisis significativo. Necesitas al menos 2 grupos.")
194             else:
195                 fig3 = px.box(df_box, x=col_region, y=
        col_resultado_comun,
196                             title=f"Distribución de {
        col_resultado_comun} por {col_region}")
197                 st.plotly_chart(fig3, use_container_width=True)
198                 st.markdown(f"""
199                     **Interpretación:** Los diagramas de caja
        muestran la distribución de **'{col_resultado_comun}'**
        para cada categoría de **'{col_region}'**.
        Puedes identificar:
        200                     - **Mediana:** La línea central de la caja.
        201

```

```

202         - **Rango intercuartílico (IQR):** La altura de
    la caja (donde se concentra el 50% central de los datos).
203         - **Bigotes:** La dispersión de los datos fuera
    del IQR.
204         - **Puntos (Outliers):** Valores atípicos que
    están significativamente fuera del rango.
205         Esto te ayuda a ver rápidamente si hay
    diferencias en la tendencia central y la dispersión entre regiones.
206         """
207
208     # --- 4. Visualización Combinada: Mapa de Calor de Medias
    ---
209     st.subheader("4          Mapa de Calor de Medias por Dos
    Categorías")
210     st.markdown("Este gráfico muestra cómo el promedio de una
    variable numérica se distribuye entre las combinaciones de dos
    variables categóricas (o discretas).")
211
212     col_heatmap_x = st.selectbox("Selecciona la columna para el
    Eje X (Categoría 1)", available_columns, key="heatmap_x")
213     col_heatmap_y = st.selectbox("Selecciona la columna para el
    Eje Y (Categoría 2)", available_columns, key="heatmap_y")
214     col_heatmap_value = st.selectbox("Selecciona la columna
    numérica para el valor (Ej. Resultado, Ingreso)", opciones_resultado,
    key="heatmap_value")
215
216     if col_heatmap_x and col_heatmap_y and col_heatmap_value:
217         if col_heatmap_x not in df.columns or col_heatmap_y not
    in df.columns or col_heatmap_value not in df.columns:
218             st.error("Una de las columnas seleccionadas para el
    mapa de calor no es válida.")
219             elif not pd.api.types.is_numeric_dtype(df[
    col_heatmap_value]):
220                 st.warning("La columna de valor para el mapa de
    calor debe ser numérica.")
221             else:
222                 # Convertir a categoría si no lo son ya, para
    asegurar un comportamiento de agrupación predecible
223                 temp_df = df[[col_heatmap_x, col_heatmap_y,
    col_heatmap_value]].copy()
224                 if not pd.api.types.is_numeric_dtype(temp_df[
    col_heatmap_x]) and not pd.api.types.is_categorical_dtype(temp_df[
    col_heatmap_x]):
225                     temp_df[col_heatmap_x] = temp_df[col_heatmap_x].
    astype('category')
226                 if not pd.api.types.is_numeric_dtype(temp_df[
    col_heatmap_y]) and not pd.api.types.is_categorical_dtype(temp_df[
    col_heatmap_y]):
227                     temp_df[col_heatmap_y] = temp_df[col_heatmap_y].
    astype('category')
228
229                 pivot_table = temp_df.pivot_table(index=
    col_heatmap_y, columns=col_heatmap_x, values=col_heatmap_value,
    aggfunc='mean')
230
231                 if not pivot_table.empty:
232                     fig_heatmap_custom = px.imshow(
    pivot_table,
233

```

```

234         text_auto=".2f",
235         aspect="auto",
236         color_continuous_scale=px.colors.sequential.
Viridis,
237         title=f"Promedio de {col_heatmap_value} por
{col_heatmap_y} y {col_heatmap_x}"
238     )
239     fig_heatmap_custom.update_xaxes(side="top")
240     fig_heatmap_custom.update_layout(height=600,
width=800)
241     st.plotly_chart(fig_heatmap_custom,
use_container_width=True)
242     st.markdown(f"""
243         **Interpretaci n:** Este mapa de calor muestra
el promedio de **'{col_heatmap_value}'**
244         para cada combinaci n nica de **'{
col_heatmap_x}'** y **'{col_heatmap_y}'**.
245         Los colores m s intensos (o m s claros,
dependiendo de la escala) indican promedios m s altos (o bajos).
246         Es til para identificar r pidamente
combinaciones de categor as con resultados particulares.
247         """)
248     else:
249         st.warning("No hay suficientes datos para
generar el mapa de calor con las columnas seleccionadas despu s de
la agrupaci n.")
250     else:
251         st.info("Selecciona tres columnas (dos categ ricas/
discretas y una num rica) para generar el mapa de calor de medias.")
252
253
254         st.markdown("         **Interpretaci n General**:: Puedes ajustar
los campos para personalizar el an lisis seg n tu CSV.")
255
256     except Exception as e:
257         st.error(f"Ocurri un error inesperado al procesar los datos o
generar gr ficos: {e}")
258         st.info("Aseg rate de que tu archivo CSV tenga el formato
esperado y que las columnas seleccionadas sean apropiadas para el
tipo de gr fico.")
259 else:
260     st.info("Por favor, sube un archivo CSV para comenzar.")

```

Listing 3: Tercera Aplicaci3n

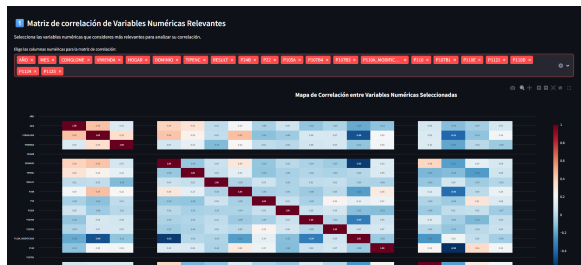
3.2. Capturas de Pantalla



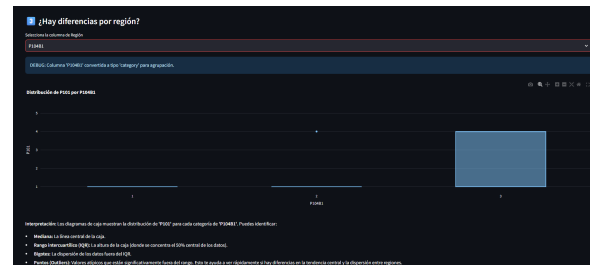
(a) [Descripción de la imagen 1]



(b) [Descripción de la imagen 2]



(a) [Descripción de la imagen 3]



(b) [Descripción de la imagen 4]

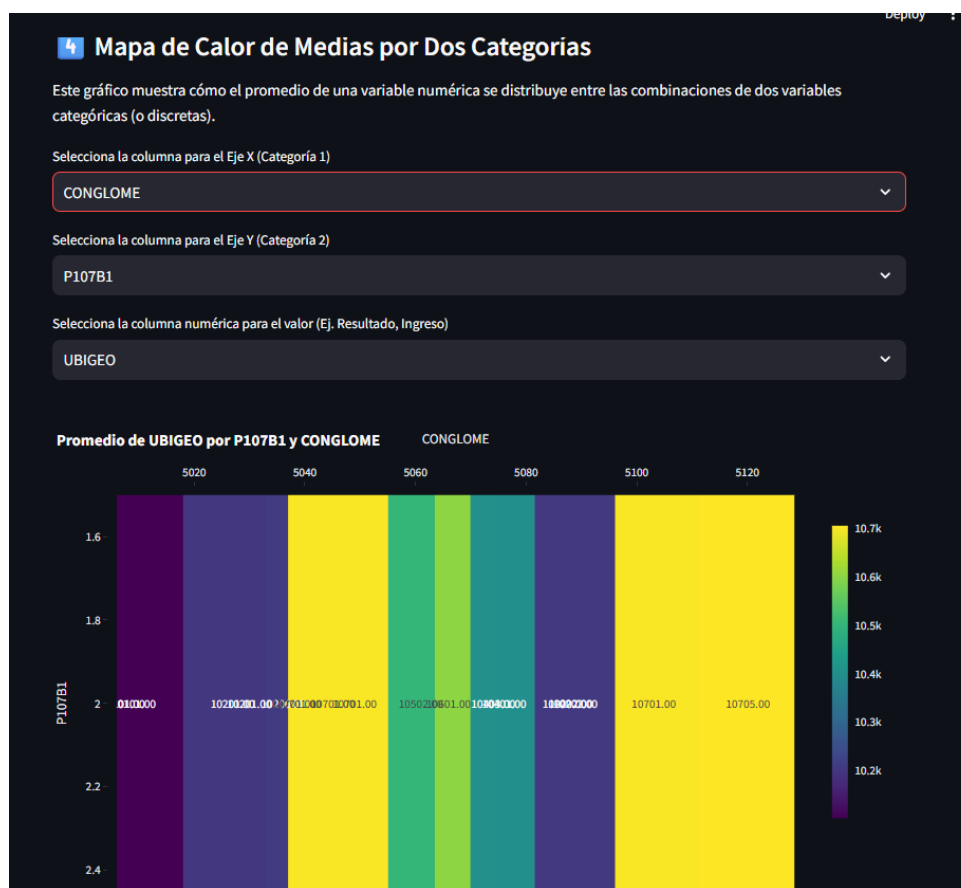


Figura 9: [Descripción de la imagen 5]

3.3. Análisis Técnico

Aspectos técnicos destacados:

- **Framework Streamlit:** Permite crear rápidamente una **interfaz web interactiva** a partir de código Python, simplificando el desarrollo de la UI.
- **Gestión de Datos con Pandas:** Utiliza **Pandas** para un manejo eficiente de DataFrames, incluyendo carga robusta (multi-codificación), limpieza y preprocesamiento de datos.
- **Visualizaciones Interactivas con Plotly:** Emplea **Plotly Express** y **Graph Objects** para generar gráficos dinámicos y explorables, mejorando la interacción del usuario con los datos.
- **Manejo Robusto de Errores y Flexibilidad:** Incluye lógica para **adaptarse a diversos CSVs** (limpieza de columnas, detección de tipos de datos) y proporciona retroalimentación clara ante errores.
- **Arquitectura Basada en Componentes:** La aplicación se construye a partir de **componentes modulares de Streamlit**, lo que facilita la adición o modificación de nuevas funcionalidades de análisis.

4. Análisis de Gastos en Servicios Básicos con MANOVA

Propósito: Permitir a los usuarios ****analizar diferencias significativas en patrones de gasto multivariados**** (varios servicios básicos a la vez) entre distintos grupos (ej. estratos socioeconómicos), ofreciendo una plataforma interactiva para la ejecución e interpretación del MANOVA.

Funcionalidades principales:

- **Carga Flexible de Datos:** Admite archivos CSV con manejo automático de múltiples codificaciones y separadores, o permite usar datos de ejemplo sintéticos.
- **Preprocesamiento Automatizado:** Realiza limpieza de nombres de columnas y validación de datos para asegurar la idoneidad para MANOVA.
- **Análisis MANOVA Interactivo:** Ejecuta el Análisis Multivariado de Varianza y presenta sus resultados clave (ej. Traza de Pillai, p-valor).
- **Visualizaciones Exploratorias y de Diagnóstico:** Ofrece gráficos (pastel, barras, violín) y métricas para entender la distribución de datos y verificar supuestos.
- **Interpretación Guiada y Post-Hoc:** Proporciona explicaciones claras de los resultados del MANOVA y realiza análisis ANOVA individuales por variable para un detalle adicional.

Características técnicas:

- **Framework Streamlit:** Facilita la creación de una **interfaz web intuitiva** y de fácil despliegue en Python.
- **Gestión de Datos con Pandas:** Utiliza **Pandas** para la manipulación, validación y preparación robusta de los datos.
- **Análisis Estadístico con Statsmodels y SciPy:** Implementa el MANOVA con ‘statsmodels.multivariate.manova’ y un método de respaldo con ‘scipy.stats’ para mayor robustez.
- **Visualización Interactiva con Plotly Express:** Genera **gráficos dinámicos** (mapas de calor, barras, violín) que permiten la exploración de los datos.
- **Manejo Inteligente de Errores:** Ofrece **mensajes contextuales** (errores, advertencias) para guiar al usuario ante problemas en los datos o el análisis.

4.1. Código Fuente

```
1 import streamlit as st
2 import pandas as pd
3 import numpy as np
4 import plotly.express as px
5 import plotly.graph_objects as go
6 from plotly.subplots import make_subplots
7 from scipy.stats import multivariate_normal
8 from scipy import stats
```

```

9 import seaborn as sns
10 import matplotlib.pyplot as plt
11 from sklearn.preprocessing import StandardScaler
12 from statsmodels.multivariate.manova import MANOVA
13 import warnings
14 warnings.filterwarnings('ignore')
15
16 # Configuraci n de la p gina
17 st.set_page_config(
18     page_title="Storytelling: An lisis MANOVA de Servicios B sicos",
19     page_icon=" ",
20     layout="wide",
21     initial_sidebar_state="expanded"
22 )
23
24 # T tulo principal
25 st.title("      Storytelling con MANOVA: An lisis de Gastos en
    Servicios B sicos")
26 st.markdown("---")
27
28 # Funci n para generar datos sint ticos (como ejemplo)
29 @st.cache_data
30 def generar_datos_ejemplo():
31     """Genera datos sint ticos como ejemplo"""
32     np.random.seed(42)
33     n_hogares = 300
34
35     # Definir grupos (estratos socioecon micos)
36     grupos = np.random.choice(['Bajo', 'Medio', 'Alto'], n_hogares, p
    =[0.4, 0.4, 0.2])
37
38     # Generar gastos base seg n grupo
39     gastos_base = {
40         'Bajo': {'agua': 50, 'luz': 80, 'gas': 40, 'internet': 30, '
    telefono': 25},
41         'Medio': {'agua': 80, 'luz': 150, 'gas': 70, 'internet': 60, '
    telefono': 45},
42         'Alto': {'agua': 120, 'luz': 250, 'gas': 100, 'internet': 100, '
    telefono': 70}
43     }
44
45     # Crear DataFrame
46     data = []
47     for i, grupo in enumerate(grupos):
48         base = gastos_base[grupo]
49         # Agregar variabilidad realista
50         agua = np.random.normal(base['agua'], base['agua'] * 0.3)
51         luz = np.random.normal(base['luz'], base['luz'] * 0.4)
52         gas = np.random.normal(base['gas'], base['gas'] * 0.35)
53         internet = np.random.normal(base['internet'], base['internet'] *
    0.2)
54         telefono = np.random.normal(base['telefono'], base['telefono'] *
    0.25)
55
56     # Asegurar valores positivos
57     agua = max(20, agua)
58     luz = max(30, luz)
59     gas = max(15, gas)

```

```

60     internet = max(10, internet)
61     telefono = max(10, telefono)
62
63     data.append({
64         'hogar_id': i+1,
65         'grupo': grupo,
66         'agua': round(agua, 2),
67         'luz': round(luz, 2),
68         'gas': round(gas, 2),
69         'internet': round(internet, 2),
70         'telefono': round(telefono, 2),
71         'total': round(agua + luz + gas + internet + telefono, 2)
72     })
73
74     return pd.DataFrame(data)
75
76 # Funci n para validar datos
77 def validar_datos(df):
78     """Valida que los datos sean apropiados para MANOVA"""
79     errores = []
80
81     if df is None or df.empty:
82         errores.append("El archivo est vac o")
83         return errores
84
85     # Verificar que tenga al menos 3 filas
86     if len(df) < 10:
87         errores.append("Se necesitan al menos 10 observaciones para un
an lisis confiable")
88
89     # Verificar columnas num ricas
90     columnas_numericas = df.select_dtypes(include=[np.number]).columns
91     if len(columnas_numericas) < 2:
92         errores.append("Se necesitan al menos 2 variables num ricas
para MANOVA")
93
94     return errores
95
96 # Funci n para limpiar nombres de columnas
97 def limpiar_nombres_columnas(df):
98     """Limpia los nombres de columnas para evitar problemas con MANOVA
"""
99     df_limpio = df.copy()
100
101     # Diccionario para mapear nombres originales a nombres limpios
102     mapeo_nombres = {}
103
104     for col in df_limpio.columns:
105         # Limpiar el nombre: solo letras, n meros y guiones bajos
106         nombre_limpio = ''.join(c if c.isalnum() or c == '_' else '_'
for c in str(col))
107         # Asegurar que empiece con letra
108         if nombre_limpio and not nombre_limpio[0].isalpha():
109             nombre_limpio = 'var_' + nombre_limpio
110         # Evitar nombres vac os
111         if not nombre_limpio:
112             nombre_limpio = f'variable_{len(mapeo_nombres)}'
113

```



```

114     mapeo_nombres[col] = nombre_limpio
115
116     # Renombrar columnas
117     df_limpio = df_limpio.rename(columns=mapeo_nombres)
118
119     return df_limpio, mapeo_nombres
120
121 # Funci n para realizar MANOVA
122 def realizar_manova(df, variables_dependientes, variable_grupo):
123     """Realiza el an lisis MANOVA"""
124     try:
125         # Limpiar nombres de columnas
126         df_limpio, mapeo = limpiar_nombres_columnas(df)
127
128         # Mapear nombres de variables a nombres limpios
129         variables_limpias = [mapeo[var] for var in
variables_dependientes]
130         grupo_limpio = mapeo[variable_grupo]
131
132         # Verificar que tengamos suficientes datos
133         if len(df_limpio) < 10:
134             return None, "Se necesitan al menos 10 observaciones para un
an lisis confiable"
135
136         # Verificar que cada grupo tenga suficientes observaciones
137         conteo_grupos = df_limpio[grupo_limpio].value_counts()
138         if conteo_grupos.min() < 3:
139             return None, f"Cada grupo debe tener al menos 3
observaciones. Grupo con menos datos: {conteo_grupos.min()}"
140
141         # Eliminar filas con valores faltantes
142         df_limpio = df_limpio.dropna(subset=variables_limpias + [
grupo_limpio])
143
144         if len(df_limpio) == 0:
145             return None, "No quedan datos despu s de eliminar valores
faltantes"
146
147         # Crear f rmula para MANOVA
148         formula = f"{' + '.join(variables_limpias)} ~ {grupo_limpio}"
149
150         # Realizar MANOVA
151         manova = MANOVA.from_formula(formula, data=df_limpio)
152         resultado = manova.mv_test()
153
154         return resultado, None
155
156     except Exception as e:
157         # Intentar m todo alternativo con scipy
158         try:
159             return realizar_manova_scipy(df, variables_dependientes,
variable_grupo)
160         except Exception as e2:
161             return None, f"Error en MANOVA: {str(e)}. Error alternativo:
{str(e2)}"
162
163 # Funci n alternativa usando scipy
164 def realizar_manova_scipy(df, variables_dependientes, variable_grupo):

```

```

165 """M todo alternativo para MANOVA usando scipy"""
166 from scipy.stats import f
167
168 # Preparar datos
169 grupos = df[variable_grupo].unique()
170 n_grupos = len(grupos)
171
172 if n_grupos < 2:
173     return None, "Se necesitan al menos 2 grupos para el an lisis"
174
175 # Calcular matrices necesarias para MANOVA
176 X = df[variables_dependientes].values
177 y = df[variable_grupo].values
178
179 # Calcular medias por grupo
180 medias_grupo = {}
181 for grupo in grupos:
182     mask = y == grupo
183     medias_grupo[grupo] = X[mask].mean(axis=0)
184
185 # Media global
186 media_global = X.mean(axis=0)
187
188 # Matriz de suma de cuadrados entre grupos (H)
189 H = np.zeros((len(variables_dependientes), len(
variables_dependientes)))
190 for grupo in grupos:
191     mask = y == grupo
192     n_grupo = mask.sum()
193     diff = medias_grupo[grupo] - media_global
194     H += n_grupo * np.outer(diff, diff)
195
196 # Matriz de suma de cuadrados dentro de grupos (E)
197 E = np.zeros((len(variables_dependientes), len(
variables_dependientes)))
198 for grupo in grupos:
199     mask = y == grupo
200     X_grupo = X[mask]
201     for i in range(len(X_grupo)):
202         diff = X_grupo[i] - medias_grupo[grupo]
203         E += np.outer(diff, diff)
204
205 # Calcular estad sticos
206 try:
207     # Pillai's Trace
208     eigenvals = np.linalg.eigvals(H @ np.linalg.inv(E))
209     pillai_trace = np.sum(eigenvals / (1 + eigenvals))
210
211     # Aproximaci n F para Pillai's Trace
212     p = len(variables_dependientes)
213     df_hyp = p * (n_grupos - 1)
214     df_err = len(df) - n_grupos
215
216     # Asegurar que df_err - p + 1 no sea cero o negativo
217     den_df_pillai = max(1, df_err - p + 1)
218
219     # Evitar divisi n por cero para el c lculo de f_stat si
pillai_trace es 1 (separaci n perfecta)

```

```

220     if (1 - pillai_trace) == 0:
221         f_stat = np.inf
222     else:
223         f_stat = (pillai_trace / (p * (n_grupos - 1))) * ((
den_df_pillai) / (1 - pillai_trace))
224
225     # Manejar casos donde f_stat podr a ser muy grande o infinito
226     if np.isinf(f_stat):
227         p_valor = 0.0 # Altamente significativo
228     elif np.isnan(f_stat):
229         p_valor = np.nan # No se puede calcular
230     else:
231         p_valor = 1 - f.cdf(f_stat, df_hyp, den_df_pillai)
232
233     # Crear DataFrame de resultados similar al de statsmodels
234     resultado = pd.DataFrame({
235         'Value': [pillai_trace],
236         'Num DF': [df_hyp],
237         'Den DF': [den_df_pillai],
238         'F Value': [f_stat],
239         'Pr > F': [p_valor]
240     }, index=[variable_grupo])
241
242     return resultado, None
243
244     except Exception as e:
245         return None, f"Error en calculo alternativo: {str(e)}"
246
247 # Sidebar para configuraci n
248 st.sidebar.header("          Configuraci n del An lisis")
249
250 # Opci n para cargar datos
251 opcion_datos = st.sidebar.radio(
252     "    Cmo    quieres cargar los datos?",
253     ["          Subir mi archivo CSV", "          Usar datos de ejemplo"]
254 )
255
256 # Cargar datos
257 df = None
258 if opcion_datos == "          Subir mi archivo CSV":
259     st.sidebar.subheader("          Carga tu archivo CSV")
260     archivo_csv = st.sidebar.file_uploader(
261         "Selecciona tu archivo CSV",
262         type=['csv'],
263         help="El archivo debe contener columnas num ricas para gastos y
una columna categ rica para agrupar"
264     )
265
266 # Opciones adicionales para la carga
267 with st.sidebar.expander("          Opciones avanzadas de carga"):
268     separador = st.selectbox(
269         "Separador del CSV:",
270         [",", ";", "\t", "|"],
271         help="Car cter que separa las columnas"
272     )
273
274     decimal = st.selectbox(
275         "Separador decimal:",

```

```

276         [".", ",",],
277         help="Carácter usado para decimales"
278     )
279
280     if archivo_csv is not None:
281         try:
282             # Intentar diferentes codificaciones
283             codificaciones = ['utf-8', 'latin-1', 'iso-8859-1', 'cp1252',
284                             , 'utf-16']
285             df = None
286             codificacion_exitosa = None
287
288             for encoding in codificaciones:
289                 try:
290                     # Resetear el puntero del archivo
291                     archivo_csv.seek(0)
292                     df = pd.read_csv(archivo_csv, encoding=encoding, sep=
293                                     =separador, decimal=decimal)
294                     codificacion_exitosa = encoding
295                     break
296                 except UnicodeDecodeError:
297                     continue
298                 except Exception as e:
299                     # Si es otro tipo de error, intentar con la
300                     siguiente codificación
301                     continue
302
303             if df is not None:
304                 st.sidebar.success(f"        Archivo cargado: {len(df)}
305                                     filas")
306                 st.sidebar.info(f"        Codificación detectada: {
307                                 codificacion_exitosa}")
308             else:
309                 st.sidebar.error("        No se pudo leer el archivo con
310                                 ninguna codificación estándar")
311                 st.sidebar.markdown("""
312                                     **Posibles soluciones:**
313                                     1. Guarda el archivo como CSV UTF-8 desde Excel
314                                     2. Usa un editor de texto para cambiar la codificación
315                                     3. Verifica que el archivo sea realmente un CSV
316                                     """)
317
318             except Exception as e:
319                 st.sidebar.error(f"        Error al cargar el archivo: {str(e)}")
320         )
321
322         st.sidebar.markdown("""
323             **Posibles causas:**
324             - Formato de archivo incorrecto
325             - Archivo corrupto
326             - Separadores no estándar
327             """)
328
329         else:
330             st.info("        Por favor, sube tu archivo CSV usando el panel
331                     lateral")
332             st.markdown("""
333                 ###        Formato requerido del CSV:
334
335                 Tu archivo debe tener:

```

```

326     - **Una columna categórica** para agrupar (ej: 'grupo', '
estrato', 'region')
327     - **Varias columnas numéricas** con los gastos en servicios (ej
: 'agua', 'luz', 'gas', etc.)
328
329     **Ejemplo de estructura:**
330     '''
331     grupo,agua,luz,gas,internet,telefono
332     Bajo,45.50,85.20,35.10,25.00,20.50
333     Medio,78.30,145.80,65.40,55.20,42.10
334     Alto,115.60,245.90,95.80,95.50,68.30
335     '''
336
337     ### Problemas comunes y soluciones:
338
339     **Si tienes problemas de codificación:**
340     - Abre tu archivo en Excel
341     - Ve a "Archivo" "Guardar como"
342     - Selecciona "CSV UTF-8 (delimitado por comas)"
343     - Guarda con ese formato
344
345     **Si tienes separadores diferentes:**
346     - Usa las opciones avanzadas en el panel lateral
347     - Prueba con punto y coma (;) si tu archivo viene de Excel en
español
348
349     **Si tienes decimales con coma:**
350     - Cambia el separador decimal en las opciones avanzadas
351     """)
352
353     # Mostrar ejemplo de archivo descargable
354     st.markdown("### Descargar archivo de ejemplo")
355     ejemplo_data = {
356         'grupo': ['Bajo', 'Bajo', 'Medio', 'Medio', 'Alto', 'Alto'],
357         'agua': [45.50, 48.20, 78.30, 82.10, 115.60, 120.80],
358         'luz': [85.20, 90.10, 145.80, 150.20, 245.90, 260.30],
359         'gas': [35.10, 38.50, 65.40, 68.20, 95.80, 100.10],
360         'internet': [25.00, 28.30, 55.20, 58.90, 95.50, 98.20],
361         'telefono': [20.50, 22.10, 42.10, 45.30, 68.30, 72.40]
362     }
363     ejemplo_df = pd.DataFrame(ejemplo_data)
364     csv_ejemplo = ejemplo_df.to_csv(index=False)
365     st.download_button(
366         label="Descargar CSV de ejemplo",
367         data=csv_ejemplo,
368         file_name="ejemplo_servicios_basicos.csv",
369         mime="text/csv"
370     )
371 else:
372     df = generar_datos_ejemplo()
373     st.sidebar.success("Usando datos de ejemplo")
374
375 # Continuar solo si hay datos
376 if df is not None:
377     # Validar datos
378     errores = validar_datos(df)
379     if errores:
380         st.error("Errores en los datos:")

```

```

381         for error in errores:
382             st.error(f"         {error}")
383         st.stop()
384
385     # Configuraci n del an lisis
386     st.sidebar.subheader("         Configuraci n del An lisis")
387
388     # Seleccionar columnas
389     columnas_disponibles = df.columns.tolist()
390     columnas_numericas = df.select_dtypes(include=[np.number]).columns.
391     tolist()
392     columnas_categoricas = df.select_dtypes(include=['object', 'category
393     ']).columns.tolist()
394
395     # Variable de agrupaci n
396     if columnas_categoricas:
397         variable_grupo = st.sidebar.selectbox(
398             "Selecciona la variable de agrupaci n:",
399             columnas_categoricas,
400             help="Variable categ rica que define los grupos a comparar"
401         )
402     else:
403         st.sidebar.error("         No se encontraron columnas categ ricas")
404         st.stop()
405
406     # Variables dependientes
407     if columnas_numericas:
408         variables_dependientes = st.sidebar.multiselect(
409             "Selecciona las variables de gasto:",
410             columnas_numericas,
411             default=columnas_numericas[:5] if len(columnas_numericas) >=
412             5 else columnas_numericas,
413             help="Variables num ricas que representan los gastos en
414             servicios"
415         )
416     else:
417         if len(variables_dependientes) < 2:
418             st.sidebar.error("         Selecciona al menos 2 variables de
419             gasto")
420             st.stop()
421         else:
422             st.sidebar.error("         No se encontraron columnas num ricas")
423             st.stop()
424
425     # Otras configuraciones
426     mostrar_datos = st.sidebar.checkbox("Mostrar datos crudos", value=
427     False)
428     nivel_confianza = st.sidebar.slider("Nivel de confianza", 0.90,
429     0.99, 0.95, 0.01)
430
431     # Filtrar datos seg n selecci n
432     df_analisis = df[[variable_grupo] + variables_dependientes].copy()
433
434     # Eliminar filas con valores faltantes
435     df_analisis = df_analisis.dropna()
436
437     # Verificar que haya al menos 2 grupos
438     grupos_unicos = df_analisis[variable_grupo].unique()

```

```

432     if len(grupos_unicos) < 2:
433         st.error("        Se necesitan al menos 2 grupos diferentes para el
an lisis")
434         st.stop()
435
436     # Crear columna total si no existe
437     if 'total' not in df_analisis.columns:
438         df_analisis['total'] = df_analisis[variables_dependientes].sum(
axis=1)
439
440     # Secci n 1: Introducci n y contexto
441     st.header("        Cap tulo 1: El Contexto del Problema")
442     st.markdown(f"""
443     ### La Historia detr s de los N meros
444
445     En este an lisis, exploraremos los patrones de gasto en servicios
b sicos entre diferentes grupos.
446     Utilizaremos **{len(df_analisis)} observaciones** divididas en **{
len(grupos_unicos)} grupos** para entender si existen diferencias
significativas en el comportamiento de gasto.
447
448     **Nuestra pregunta de investigaci n:** Existen diferencias
significativas en los patrones de gasto
449     en servicios b sicos entre los diferentes grupos analizados?
450
451     **Variables analizadas:**
452     - **Variable de agrupaci n:** {variable_grupo}
453     - **Variables de gasto:** {', '.join(variables_dependientes)}
454
455     Para responder esta pregunta, utilizaremos **MANOVA (An lisis
Multivariado de Varianza)** ,
456     una t cnica estad stica que nos permite analizar m ltiples
variables dependientes simult neamente.
457     """)
458
459     if mostrar_datos:
460         st.subheader("        Datos del An lisis")
461         st.dataframe(df_analisis)
462
463         # Estad sticas descriptivas
464         st.subheader("        Estad sticas Descriptivas")
465         stats_desc = df_analisis.groupby(variable_grupo)[
variables_dependientes].agg(['count', 'mean', 'std', 'min', 'max']).
round(2)
466         st.dataframe(stats_desc)
467
468     # Secci n 2: An lisis Exploratorio
469     st.header("        Cap tulo 2: Explorando los Patrones de Gasto")
470
471     col1, col2 = st.columns(2)
472
473     with col1:
474         st.subheader("Distribuci n por Grupo")
475         conteo_grupos = df_analisis[variable_grupo].value_counts()
476         fig_dist = px.pie(values=conteo_grupos.values, names=
conteo_grupos.index,
477                             title=f"Distribuci n de Observaciones por {
variable_grupo}")

```

```

478     st.plotly_chart(fig_dist, use_container_width=True)
479
480     with col2:
481         st.subheader("Gasto Total Promedio por Grupo")
482         gasto_promedio = df_analisis.groupby(variable_grupo)['total'].
mean().reset_index()
483         fig_bar = px.bar(gasto_promedio, x=variable_grupo, y='total',
484                         title="Gasto Total Promedio por Grupo",
485                         color=variable_grupo)
486         st.plotly_chart(fig_bar, use_container_width=True)
487
488     # An lisis por servicios individuales
489     st.subheader("          An lisis Detallado por Variable")
490
491     servicio_seleccionado = st.selectbox("Selecciona una variable para
analizar:", variables_dependientes)
492
493     fig_violin = px.violin(df_analisis, x=variable_grupo, y=
servicio_seleccionado,
494                          title=f"Distribuci n de {
servicio_seleccionado}",
495                          color=variable_grupo)
496     st.plotly_chart(fig_violin, use_container_width=True)
497
498     # Informaci n de diagn stico
499     st.subheader("          Diagn stico de Datos")
500
501     col1, col2, col3 = st.columns(3)
502
503     with col1:
504         st.metric("Total de observaciones", len(df_analisis))
505         st.metric("Observaciones despu s de limpiar", len(df_analisis.
dropna()))
506
507     with col2:
508         st.metric("N mero de grupos", len(grupos_unicos))
509         min_obs_grupo = df_analisis[variable_grupo].value_counts().min()
510         st.metric("M nimo obs. por grupo", min_obs_grupo)
511
512     with col3:
513         st.metric("Variables analizadas", len(variables_dependientes))
514         valores_faltantes = df_analisis[variables_dependientes].isnull()
.sum().sum()
515         st.metric("Valores faltantes", valores_faltantes)
516
517     # Mostrar distribuci n por grupo
518     st.subheader("          Distribuci n de Observaciones por Grupo")
519     dist_grupos = df_analisis[variable_grupo].value_counts().sort_index
()
520     st.dataframe(dist_grupos.to_frame("Cantidad"))
521
522     # Advertencias
523     if min_obs_grupo < 3:
524         st.warning(f"          Algunos grupos tienen menos de 3
observaciones. M nimo: {min_obs_grupo}")
525
526     if valores_faltantes > 0:
527         st.warning(f"          Hay {valores_faltantes} valores faltantes

```



```

que ser n eliminados del an lisis")
528
529 # Verificar variabilidad
530 variabilidad = df_analisis[variables_dependientes].std()
531 variables_sin_variabilidad = variabilidad[variabilidad == 0].index.
tolist()
532 if variables_sin_variabilidad:
533     st.warning(f"          Variables sin variabilidad: {'', '.join(
variables_sin_variabilidad)}")
534
535 # Mostrar primeras filas para verificar datos
536 with st.expander("          Vista previa de datos"):
537     st.dataframe(df_analisis.head(10))
538
539 # Secci n 3: An lisis MANOVA
540 st.header("          Cap tulo 3: El An lisis MANOVA")
541
542 st.markdown("""
543 ###    Qu     es MANOVA?
544
545 **MANOVA (Multivariate Analysis of Variance)** es una extensi n del
ANOVA que nos permite:
546 - Analizar m ltiples variables dependientes simult neamente
547 - Controlar el error tipo I al hacer m ltiples comparaciones
548 - Detectar diferencias que podr an no ser evidentes en an lisis
univariados
549
550 **Hip tesis de nuestro estudio:**
551 - ** H  :** No hay diferencias significativas en los gastos entre
grupos
552 - ** H  :** Existen diferencias significativas en al menos una
variable entre grupos
553 """)
554
555 # Realizar MANOVA
556 with st.spinner("Realizando an lisis MANOVA..."):
557     resultado_manova, error_manova = realizar_manova(df_analisis,
variables_dependientes, variable_grupo)
558
559 if error_manova:
560     st.error(f"          Error en el an lisis MANOVA: {error_manova}")
561     st.markdown("""
562     **Posibles causas del error:**
563     - Nombres de columnas con caracteres especiales
564     - Muy pocas observaciones por grupo (m nimo 3 por grupo)
565     - Variables con poca variabilidad
566     - Problemas de multicolinealidad
567     - Valores faltantes en los datos
568
569     **Recomendaciones:**
570     1. Verifica que cada grupo tenga al menos 3 observaciones
571     2. Revisa que no haya valores faltantes
572     3. Aseg rate de que las variables num ricas tengan
variabilidad
573     4. Simplifica los nombres de columnas (sin espacios ni
caracteres especiales)
574     """)
575 else:

```

```

576         st.subheader("Resultados del An lisis MANOVA")
577
578         # Mostrar resultados
579         st.markdown("### Estad sticos de Prueba")
580         # Acceder a los DataFrames de los resultados (Pillai's Trace,
581         # Wilks' Lambda, etc.) desde el diccionario de resultados.
582         # El m todo 'mv_test()' usualmente retorna un diccionario con
583         # claves como 'Target', 'Intercept', etc.
584         # Estamos interesados en la tabla para la variable independiente
585         (variable_grupo).
586
587         # Esto obtendr el DataFrame para tu variable independiente
588         if isinstance(resultado_manova, dict) and variable_grupo in
589         resultado_manova:
590             resultado_df = resultado_manova[variable_grupo].round(4)
591         elif isinstance(resultado_manova, pd.DataFrame): # Para el caso
592         de fallback de scipy
593             resultado_df = resultado_manova.round(4)
594         else:
595             st.warning("No se pudo formatear los resultados del MANOVA.
596             Mostrando el objeto completo.")
597             st.write(resultado_manova)
598             resultado_df = pd.DataFrame() # Crear un DataFrame vac o
599             para evitar errores posteriores
600
601         if not resultado_df.empty:
602             st.dataframe(resultado_df)
603
604         # Interpretaci n de resultados
605         st.markdown("### Interpretaci n de los Resultados")
606
607         # Obtener el p-valor (usando diferentes m todos seg n el
608         resultado)
609         # Priorizar 'Pr > F' y manejar casos donde la estructura
610         podr a ser diferente
611         p_valor_pillai = np.nan
612         if 'Pr > F' in resultado_df.columns:
613             p_valor_pillai = resultado_df.loc[resultado_df.index[0],
614             "Pr > F"]
615         elif 'P-value' in resultado_df.columns: # Para posibles
616         cambios futuros u otras librer as
617             p_valor_pillai = resultado_df.loc[resultado_df.index[0],
618             "P-value"]
619         elif resultado_df.shape[1] > 0: # Intentar obtener la
620         ltima columna si no se encuentra 'Pr > F'
621             p_valor_pillai = resultado_df.iloc[0, -1]
622
623         # Asegurarse de que sea un float para la comparaci n
624         if isinstance(p_valor_pillai, (pd.Series, pd.DataFrame)):
625             p_valor_pillai = p_valor_pillai.iloc[0]
626         try:
627             p_valor_pillai = float(p_valor_pillai)
628         except (ValueError, TypeError):
629             p_valor_pillai = np.nan # Si la conversi n falla,
630         establecer a NaN
631
632         if not np.isnan(p_valor_pillai) and p_valor_pillai < (1 -
633         nivel_confianza):

```

```

619         st.success(f"""
620             **      RESULTADO SIGNIFICATIVO** (p = {p_valor_pillai:.4f
        })

621
622         Con un nivel de confianza del {nivel_confianza*100}%,
        podemos **rechazar la hip tesis nula**.

623
624         **Conclusi n:** Existen diferencias estad sticamente
        significativas en los patrones de gasto
        entre los diferentes grupos analizados.
625         """
626     elif not np.isnan(p_valor_pillai):
627         st.warning(f"""
628             **      RESULTADO NO SIGNIFICATIVO** (p = {p_valor_pillai
629             :.4f})

630
631         Con un nivel de confianza del {nivel_confianza*100}%, **
        no podemos rechazar la hip tesis nula**.

632
633         **Conclusi n:** No hay evidencia suficiente para
        afirmar que existen diferencias significativas
634         en los gastos entre los grupos analizados.
635         """
636     else:
637         st.warning("No se pudo obtener un p-valor para la
        interpretaci n. Revise los resultados completos.")

638
639     # Secci n 4: An lisis Post-Hoc
640     st.header("          Cap tulo 4: An lisis Detallado por Variable")
641
642     st.markdown("""
643     Aunque MANOVA nos da una respuesta global, es importante entender **
        qu  variables espec ficas** muestran las mayores diferencias entre
        grupos.
644     """)

645
646     # An lisis univariado para cada variable
647     st.subheader("An lisis ANOVA Individual por Variable")

648
649     resultados_anova = []
650     for variable in variables_dependientes:
651         try:
652             # ANOVA para cada variable
653             grupos_variable = [df_analisis[df_analisis[variable_grupo]
        == grupo][variable].dropna()
654                             for grupo in grupos_unicos]
655             # Filtrar grupos vac os
656             grupos_variable = [grupo for grupo in grupos_variable if len
        (grupo) > 0]

657
658             # Asegurarse de que haya al menos 2 grupos no vac os y m s
        de 1 observaci n por grupo para un ANOVA significativo
659             if len(grupos_variable) >= 2 and all(len(g) > 1 for g in
        grupos_variable):
660                 f_stat, p_val = stats.f_oneway(*grupos_variable)

661
662                 resultados_anova.append({
663                     'Variable': variable,

```

```

664         'F-estadístico': round(f_stat, 4),
665         'p-valor': round(p_val, 4),
666         'Significativo': 'S' if p_val < (1 -
nivel_confianza) else 'No'
667     })
668     else:
669         resultados_anova.append({
670             'Variable': variable,
671             'F-estadístico': 'N/A',
672             'p-valor': 'N/A',
673             'Significativo': 'Datos insuficientes'
674         })
675     except Exception as e:
676         resultados_anova.append({
677             'Variable': variable,
678             'F-estadístico': 'Error',
679             'p-valor': 'Error',
680             'Significativo': 'Error'
681         })
682
683     df_anova = pd.DataFrame(resultados_anova)
684     st.dataframe(df_anova)
685
686     # Visualización de diferencias entre grupos
687     st.subheader("Comparación Visual entre Grupos")
688
689     # Crear gráfico de barras agrupadas
690     n_vars = len(variables_dependientes)
691     cols = 3
692     rows = (n_vars + cols - 1) // cols
693
694     fig_comparacion = make_subplots(
695         rows=rows, cols=cols,
696         subplot_titles=variables_dependientes,
697         vertical_spacing=0.12,
698         horizontal_spacing=0.08
699     )
700
701     for i, variable in enumerate(variables_dependientes):
702         row = i // cols + 1
703         col = i % cols + 1
704
705         medias = df_analisis.groupby(variable_grupo)[variable].mean()
706
707         fig_comparacion.add_trace(
708             go.Bar(x=medias.index, y=medias.values,
709                 name=variable,
710                 showlegend=False),
711             row=row, col=col
712         )
713
714     fig_comparacion.update_layout(height=200*rows, title_text="Gasto
Promedio por Variable y Grupo")
715     st.plotly_chart(fig_comparacion, use_container_width=True)
716
717     # Sección 5: Conclusiones
718     st.header("Capítulo 5: Conclusiones y Recomendaciones")
719

```

```

720 st.markdown(f"""
721     ###           Resumen Ejecutivo
722
723     Basado en nuestro an lisis MANOVA de **{len(df_analisis)}
observaciones**:
```

724

```

725     **Principales Hallazgos:**
726     """)
727
728     # Calcular estad sticas descriptivas por grupo
729     cols_metricas = st.columns(len(grupos_unicos))
730
731     for i, grupo in enumerate(grupos_unicos):
732         datos_grupo = df_analisis[df_analisis[variable_grupo] == grupo]
733         with cols_metricas[i]:
734             st.metric(f"Grupo: {grupo}", len(datos_grupo))
735             st.metric(f"Gasto Promedio", f"${datos_grupo['total'].mean()
:.2f}")
736             st.metric(f"Desviaci n Est ndar", f"${datos_grupo['total
'].std():.2f}")
737
738     # Recomendaciones personalizadas
739     st.markdown("""
740     ###           Recomendaciones
741
742     **Basado en el an lisis de tus datos:**
743
744     1. **Si encontraste diferencias significativas:**
745         - Investiga qu factores causan estas diferencias
746         - Considera estrategias diferenciadas por grupo
747         - Desarrolla pol ticas segmentadas para cada grupo
748
749     2. **Si no encontraste diferencias significativas:**
750         - Los grupos muestran patrones similares de gasto
751         - Puedes aplicar estrategias uniformes
752         - Considera otros factores de segmentaci n
753
754     3. **Para an lisis futuros:**
755         - Incluye m s variables explicativas
756         - Considera an lisis de series temporales
757         - Eval a factores externos que puedan influir
758     """)
759
760     # Exportar resultados
761     st.subheader("           Exportar Resultados")
762
763     col1, col2 = st.columns(2)
764
765     with col1:
766         if st.button("           Descargar Estad sticas Descriptivas"):
767             stats_export = df_analisis.groupby(variable_grupo)[
variables_dependientes].agg(['count', 'mean', 'std']).round(2)
768             csv = stats_export.to_csv()
769             st.download_button(
770                 label="           Descargar CSV",
771                 data=csv,
772                 file_name="estadisticas_descriptivas.csv",
773                 mime="text/csv"
```

```

774         )
775
776     with col2:
777         if st.button("Descargar Resultados ANOVA"):
778             csv_anova = df_anova.to_csv(index=False)
779             st.download_button(
780                 label="Descargar CSV",
781                 data=csv_anova,
782                 file_name="resultados_anova.csv",
783                 mime="text/csv"
784             )
785
786     # Sección técnica
787     with st.expander("Detalles Técnicos del Análisis"):
788         st.markdown(f"""
789         **Metodología Empleada:**
790         - **Técnica:** MANOVA (Multivariate Analysis of Variance)
791         - **Variables Dependientes:** {', '.join(variables_dependientes)}
792
793         - **Variable Independiente:** {variable_grupo}
794         - **Número de grupos:** {len(grupos_unicos)}
795         - **Tamaño de muestra:** {len(df_analisis)}
796         - **Estadístico Principal:** Pillai's Trace (más robusto ante violaciones de supuestos)
797         - **Software:** Python con librerías statsmodels, scipy, plotly
798
799         **Supuestos del MANOVA:**
800         - Normalidad multivariada
801         - Homogeneidad de matrices de covarianza
802         - Independencia de observaciones
803         - Ausencia de outliers extremos
804
805         **Grupos analizados:** {', '.join(grupos_unicos)}
806         """)
807
808     # Footer
809     st.markdown("----")
810     st.markdown("*Desarrollado con usando Python y Streamlit para análisis estadístico*")

```

Listing 4: Cuarta Aplicación

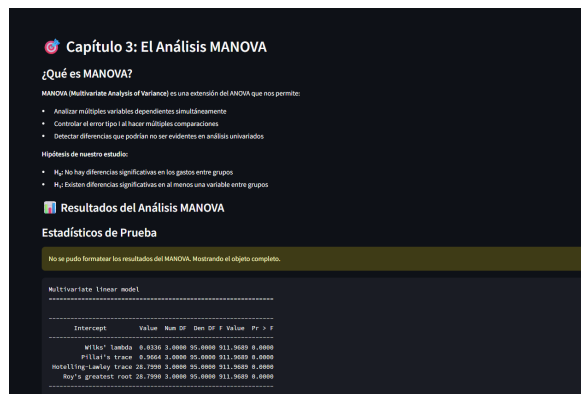
4.2. Capturas de Pantalla



(a) [Descripción de la imagen 1]



(b) [Descripción de la imagen 2]



(a) [Descripción de la imagen 3]



(b) [Descripción de la imagen 4]

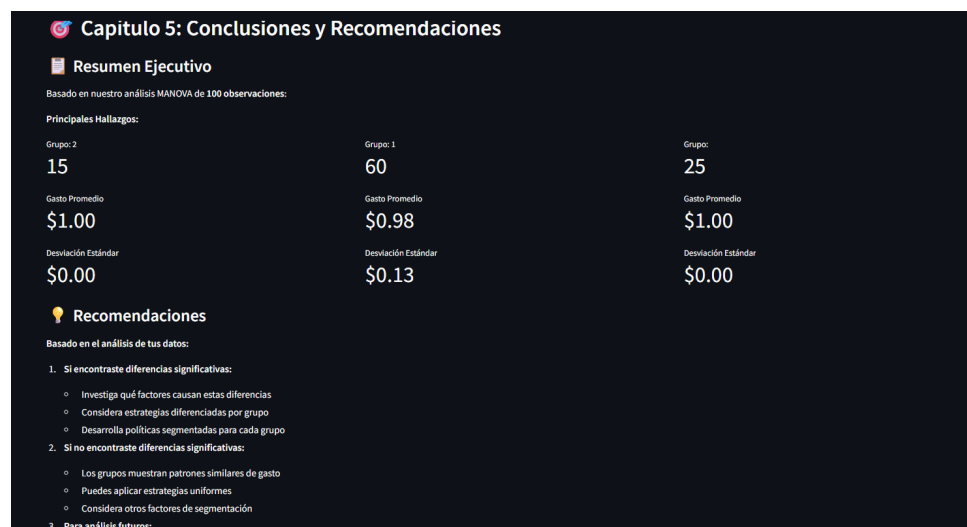


Figura 12: [Descripción de la imagen 5]

4.3. Análisis Técnico

Aspectos técnicos destacados:

- **Framework Streamlit:** Permite crear rápidamente una **interfaz web interactiva** a partir de código Python, simplificando el desarrollo de la UI.
- **Gestión de Datos con Pandas:** Utiliza **Pandas** para un manejo eficiente de DataFrames, incluyendo carga robusta (multi-codificación), limpieza y preprocesamiento de datos.
- **Análisis Estadístico con Statsmodels y SciPy:** Emplea **statsmodels** para MANOVA y **scipy.stats** para ANOVA, asegurando robustez y flexibilidad en los cálculos.
- **Visualizaciones Interactivas con Plotly:** Genera **gráficos dinámicos** y explorables (barras, violín, pastel) que enriquecen la comprensión de los datos.
- **Manejo Robusto de Errores y Adaptabilidad:** Incluye validaciones y mensajes claros para guiar al usuario, además de adaptarse a diversas estructuras de CSV.