

Visualización de evolución en el tiempo

Daniel Giovanni Rodriguez Coral

Guillermo Augusto Duran Boneth

Juan Sebastián Torres Perdomo

Fundación Universitaria Internacional de La Rioja

Especialización en Inteligencia Artificial

Ingeniero Javier Diaz Diaz

Colombia

Diciembre, 2024

Contenido

Introducción	3
Objetivos	4
Refinación de los datos	5
Visualización en Power Bi	10
Conclusiones	15

Introducción

En la actualidad, la visualización de los datos representa un papel fundamental a la hora de querer transmitir información y que ésta genere un valor agregado para el público objetivo. Desde la toma de decisiones a nivel empresarial, como también en diferentes análisis en contextos como el científico, la capacidad de poder transformar y procesar los datos en visualizaciones sumamente claras y que representen adecuadamente el objetivo que se quiere transmitir ha demostrado ser esencial.

Para alcanzar el proceso de visualización de datos, es necesario establecer de antemano una serie de pasos o criterios que faciliten una correcta manipulación de los datos. Entre estos pasos, se destacan la normalización y el procesamiento de los datos.

Otro de los puntos más importantes es definir adecuadamente las herramientas con las que se llevará a cabo el procesamiento de los datos y las visualizaciones, en este sentido, las herramientas utilizadas en el presente trabajo para el procesamiento de los datos son Python y OpenRefine, y para la construcción de las visualizaciones PowerBI.

En este documento se explora de forma detallada e integral el proceso de transformación de datos, desde la limpieza inicial hasta la visualización final, y se destaca como es que estas herramientas se complementan para poder generar valor agregado en los resultados.

Al presentar también un enfoque simple, bajo un ejemplo y resultados, este trabajo pretende evidenciar el impacto que puede llegar a tener el análisis y presentación de los datos en los contextos en que sea pertinentes aplicarlos.

Objetivos

Objetivo General

Desarrollar un dashboard interactivo para la visualización y análisis de datos, aplicando una limpieza, transformación y representación gráfica, con la finalidad de demostrar habilidades en el uso de herramientas como Power BI, Python y OpenRefine.

Objetivos Específicos

Limpiar y transformar los datos de ventas ficticias para garantizar la consistencia y precisión de la información.

Crear una visualización interactiva que permita identificar algunos hallazgos o patrones importantes dentro del conjunto de datos estudiado.

Refinación de los datos

Refinación en Python

Se realiza una limpieza previa del conjunto de datos suministrado, usando el lenguaje Python, a continuación, algunas evidencias del trabajo realizado:

1. Se crea una función que permite formatear el número de teléfono del cliente en formato Prefijo + Número de Teléfono + Extensión.

```
C:\> 2024 > Especialización IA > Visualización > Actividad_1_Visualización > Actividad_1 > Actividad_1.py > ...
1 import pandas as pd
2 import re
3
4
5 def format_telephone_number(string):
6     match = re.search(
7         r'^(?!\d{1,3})?\s\-\.\]?(\d{3})\)?\s\-\.\.]?(\d{3})[\s\-\.\.]?(\d{4})\)?\.\.?(ext\.\.?[x|#)?\s*(\d+))?',
8         string,
9         re.IGNORECASE
10    )
11    if match:
12        number = '-'.join(filter(None, match.groups()[1:4]))
13
14        prefix = match.group(1)
15        if prefix:
16            prefix = prefix.lstrip("+")
17            prefix = prefix.lstrip("0")
18            prefix = "+" + prefix
19            number = f"{prefix} {number}"
20
21        extension = match.group(6) if match.group(6) else None
22
23        if extension:
24            number = f"{number} EXT {extension}"
```

2. Se realiza la limpieza de caracteres extraños en todos los campos haciendo uso de expresiones regulares, se eliminan espacios al comienzo y al final de la palabra. También se realiza el ajuste correspondiente para el campo Fecha_Venta para que el formato sea YYYY-MM-DD

```
def format_telephone_number(string):
    return number
    return None

ruta_del_archivo = './Initial_load/datos_ventas.csv'

df = pd.read_csv(ruta_del_archivo)

df['ID_Venta'] = df['ID_Venta'].astype(str).strip()

df['Producto'] = df['Producto'].astype(str).apply(lambda text: re.sub(r'^a-zA-Z0-9 ', '', text)).str.upper()
df['Ciudad'] = df['Ciudad'].astype(str).apply(lambda text: re.sub(r'^a-zA-Z0-9 ', '', text)).str.upper()
df['Categoria'] = df['Categoria'].astype(str).apply(lambda text: re.sub(r'^a-zA-Z0-9 ', '', text)).str.upper()

df['Precio_Unitario'] = pd.to_numeric(df['Precio_Unitario'], errors='coerce')
df['Cantidad'] = pd.to_numeric(df['Cantidad'], errors='coerce')

df['Fecha_Venta'] = df['Fecha_Venta'].str.replace(r'^[^\w_]+|^[^\w_]+$', '', regex=True)
df['Fecha_Venta'] = df['Fecha_Venta'].str.replace(r'^_+|_+$', '', regex=True)
df['Fecha_Venta'] = pd.to_datetime(df['Fecha_Venta'], format='%Y-%m-%d').dt.date
```

```

df['Email'] = df['Email'].str.replace(r'^^[^w_]+|^[^w_]+$', '', regex=True)
df['Email'] = df['Email'].str.replace(r'^_+|_+$', '', regex=True)

df[['Telefono_Formateado']] = df['Telefono'].apply(lambda x: pd.Series(format_telephone_number(x)))

df['Direccion_Formateada'] = df['Direccion'].str.replace('\n', ' ', regex=False) ##La dirección tenía saldo de línea

df['Metodo_Pago'] = df['Metodo_Pago'].astype(str).apply(lambda text: re.sub(r'^a-zA-Z0-9 ', '', text)).str.upper().s

df['Estado'] = df['Estado'].astype(str).apply(lambda text: re.sub(r'^a-zA-Z0-9 ', '', text)).str.upper().str.strip()

df['Comentario'] = df['Comentario'].astype(str).apply(lambda text: re.sub(r'^a-zA-Z0-9 ', '', text)).str.upper().str

df['Descuento'] = pd.to_numeric(df['Descuento'], errors='coerce')

df['Total'] = df['Precio_Unitario'] * df['Cantidad']

df['total_con_descuento'] = df['Total'] * (1 - df['Descuento'] / 100)

print(df)

df.to_csv("./Result/datos_ventas_clean.csv", index= False)

```

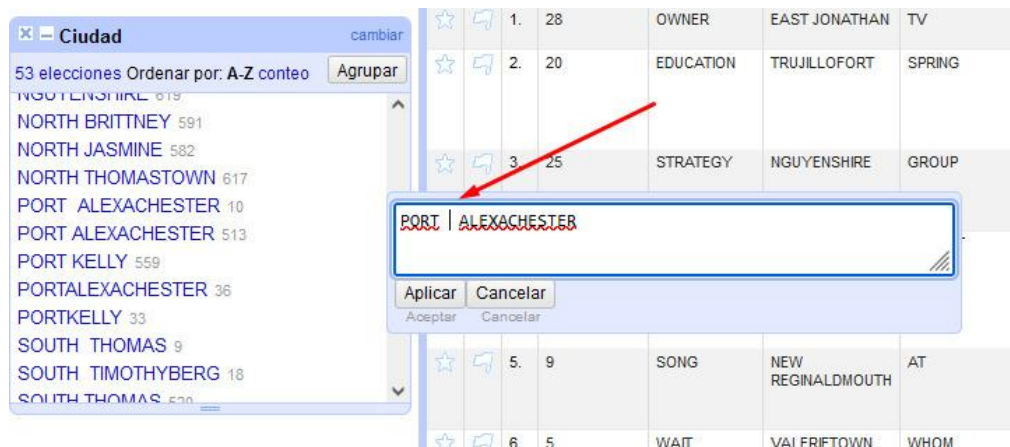
Refinación en OpenRefine

Se continúa con la organización, limpieza y creación de columnas de faltantes del conjunto de datos en OpenRefine, a continuación, algunas evidencias del trabajo realizado:

1. Se hace una revisión del formato de los emails, se verifica la falta del símbolo de arroba @, en algunos de los registros, así como espacios en blanco y caracteres que no corresponden.



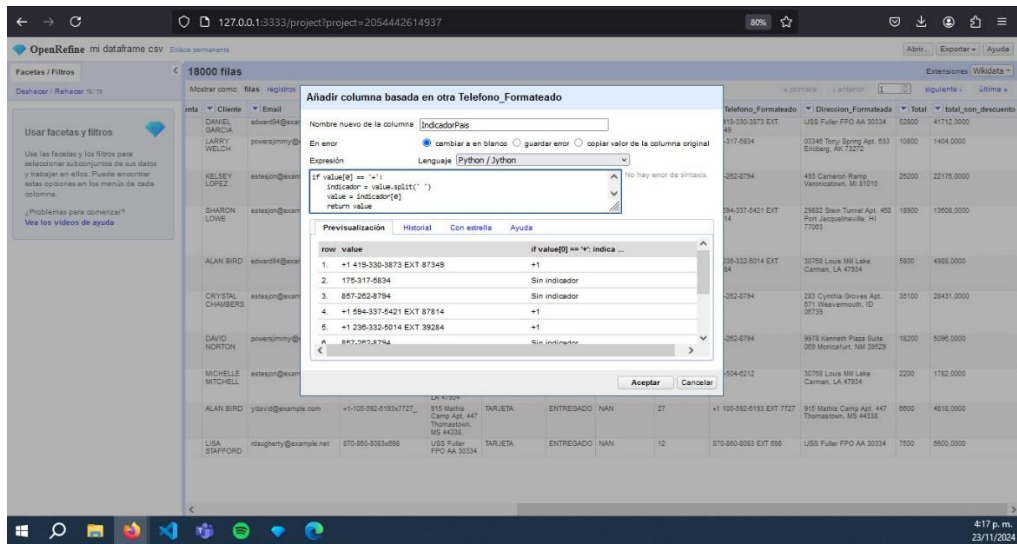
- Se estandariza los nombres de las ciudades, tomando un nombre como referencia e igualando los similares.



- Se estandariza los nombres de teléfono, eliminando caracteres que no pertenecen a dichos datos y separando los indicativos de los número telefónicos, usando código del lenguaje Python.

Configurar opciones del análisis sintáctico		Nombre del proyecto		mi dataframe csv		Etiquetas	
Id	Precio_Unitario	Cantidad	Fecha_Venta	Cliente	Email	Telefono	Direccion
	2400	22	2023-11-20	DANIEL GARCIA	edward94@example.org	001-419-330-3873x873497	USS Fuller FPO AA 30334
IG	600	18	2023-01-04	LARRY WELCH	powersjimmy@example.org	(175)317-5834	03346 Tony Spring Apt. 633 Ericberg, AK 73272
IP	900	28	2023-05-05	KELSEY LOPEZ	estesjon@example.org	857-262-8794	_493 Cameron Ramp Veroniatown, MI 81010
ET	700	27	2023-06-28	SHARON LOWE	estesjon@example.org	?001-594-337-5421x87814	29882 Stein Tunnel Apt. 468 Port Jacquelineville, 11 77062

Los datos como Codificación de caracteres US-ASCII



Añadir columna basada en otra Telefono_Formateado

Nombre nuevo de la columna

En error ☒ cambiar a en blanco ☐ guardar error ☐ copiar valor de la columna original

Expresión

```
value = data_ext[-1]
return value
else:
return 'Sin extension'
```

No hay error de sintaxis.

Previsualización Historial Con estrella Ayuda

row	value	if 'EXT' in value: data_ex ...
1.	+1 419-330-3873 EXT 87349	87349
2.	175-317-5834	Sin extension
3.	857-262-8794	Sin extension
4.	+1 594-337-5421 EXT 87814	87814
5.	+1 236-332-5014 EXT 39284	39284
6.	857-262-8794	Sin extension

Aceptar

Cancelar

5. Se estandariza los datos de la columna Cliente, transformando los nombres similares, que solamente se diferencien por cosas mínimas, por ejemplo, espacios o mayúsculas y minúsculas.

Agrupar y editar la columna "Cliente"

Busca grupos de valores de celda diferentes que puedan ser otras representaciones de la misma cosa. Por ejemplo, "Nueva York" y "nueva york" probablemente se refieran al mismo concepto y sólo se diferencien por las mayúsculas, y "Gödel" y "Godel" probablemente se refieran a la misma persona. [Conozca más...](#)

Método: **Vecino más cercano** Función distancia: **Levenshtein (distancia)** Radio: **1.0** ☐ Actualización automática **57 clusters encontrado**
Caracteres del bloque: **6**

Tamaño del grupo	Número de filas	Valores en agrupación	¿Unir?	Nuevo valor de la celda
3	618	<ul style="list-style-type: none">CHARLES WILLIAMS (565 filas)CHARLESWILLIAMS (41 filas)CHARLES WILLIAMS (12 filas)	<input type="checkbox"/>	CHARLES WILLIAMS
3	580	<ul style="list-style-type: none">JENNIFER HOLMES (540 filas)JENNIFERHOLMES (34 filas)JENNIFER HOLMES (6 filas)	<input type="checkbox"/>	JENNIFER HOLMES
3	634	<ul style="list-style-type: none">SAMANTHA HODGES (590 filas)SAMANTHAHODGES (33 filas)SAMANTHA HODGES (11 filas)	<input type="checkbox"/>	SAMANTHA HODGES
3	614	<ul style="list-style-type: none">CYNTHIA HAYES (573 filas)CYNTHIAHAYES (33 filas)CYNTHIA HAYES (8 filas)	<input type="checkbox"/>	CYNTHIA HAYES

Valores en agrupación: 2 — 3

Filas en el grupo: 520 — 640

Longitud promedio de los valores: 9.5 — 17.5

Varianza de los valores: 0.5 — 0.8170000000000001

Visualización en Power Bi

Junto a este documento se entrega el archivo con el conjunto de datos final y el dashboard realizado con los datos limpios y organizados en PowerBi. A continuación, una captura de pantalla del dashboard y la explicación de sus elementos:



1. Filtros

- **Fecha Venta:** Permite al usuario seleccionar la fecha que quiere validar, teniendo en cuenta siempre la última fecha de venta registrada en la fuente de datos. Por defecto, el dashboard siempre mostrará la información a la última fecha registrada de acuerdo al filtro aplicado por el usuario. Por ejemplo, si el usuario selecciona en el filtro Fecha Venta la fecha 30/11/2023, si la última fecha de venta registrada fue 20/11/2023, el reporte automáticamente mostrará la información hasta esta fecha.

- **Categoría:** Permite al usuario seleccionar la categoría de los productos.

- Un KPI, mostrando la última fecha de compra registrada en la fuente de datos, se incluyó para que el usuario pueda analizar la información de manera más sencilla. Siempre mostrará la última fecha de venta de acuerdo al filtro aplicado por el usuario.

- | Reporte de Ventas - Datos de contacto clientes | | | | | | |
|---|------------|--|---------------|--------------|-------------------------------|-----------------------------|
| <div> <div>Regresar al Reporte de Ventas</div> </div> | | | | | | |
| <div> <div>Cliente</div> <div>Todas</div> </div> | | <div> <div>Ciudad</div> <div>Todas</div> </div> | | | <div> <div>Reset</div> </div> | |
| CLIENTE | CIUDAD | DIRECCION | INDICADOR | TELÉFONO | EXTENSION | EMAIL |
| ADAM MUELLER | ANDREWSIDE | 025 BILLY ROW LISASIDE, AR 96338 | +1 | 557-543-1910 | SIN EXTENSION | ydauid@example.com |
| ADAM MUELLER | ANDREWSIDE | 03346 TONY SPRING APT. 633 ERICBERG, AK 73272 | SIN INDICADOR | 339-100-4282 | 04920 | rbrown@example.org |
| ADAM MUELLER | ANDREWSIDE | 03346 TONY SPRING APT. 633 ERICBERG, AK 73272 | SIN INDICADOR | 612-821-1834 | 900 | krystalmcDaniel@example.org |
| ADAM MUELLER | ANDREWSIDE | 0827 JILL CAPE SUITE 299 ANGELOTT, DC 15542 | SIN INDICADOR | 630-469-6433 | SIN EXTENSION | karaware@example.net |
| ADAM MUELLER | ANDREWSIDE | 2644 GONZALEZ MILLS APT. 856 PORT FRANK, IN 25816 | +1 | 096-236-8217 | SIN EXTENSION | brianboyd@example.org |
| ADAM MUELLER | ANDREWSIDE | 2644 GONZALEZ MILLS APT. 856 PORT FRANK, IN 25816 | SIN INDICADOR | 491-474-2586 | SIN EXTENSION | estesjon@example.org |
| ADAM MUELLER | ANDREWSIDE | 29340 MURPHY PARKWAY NORTH ROBERT, KS 42070 | SIN INDICADOR | 602-238-8997 | 526 | hornebeth@example.com |
| ADAM MUELLER | ANDREWSIDE | 29882 STEIN TUNNEL APT. 468 PORT JACQUELINEVILLE, HI 77063 | SIN INDICADOR | 369-774-6572 | 373 | |
| ADAM MUELLER | ANDREWSIDE | 30768 LOUIS MILL LAKE CARMEN, LA 47934 | +1 | 890-231-3672 | 1535 | cwilliams@example.com |
| ADAM MUELLER | ANDREWSIDE | 493 CAMERON RAMP VERONICATOWN, MI 81010 | SIN INDICADOR | 011-988-9703 | 8562 | hornebeth@example.com |

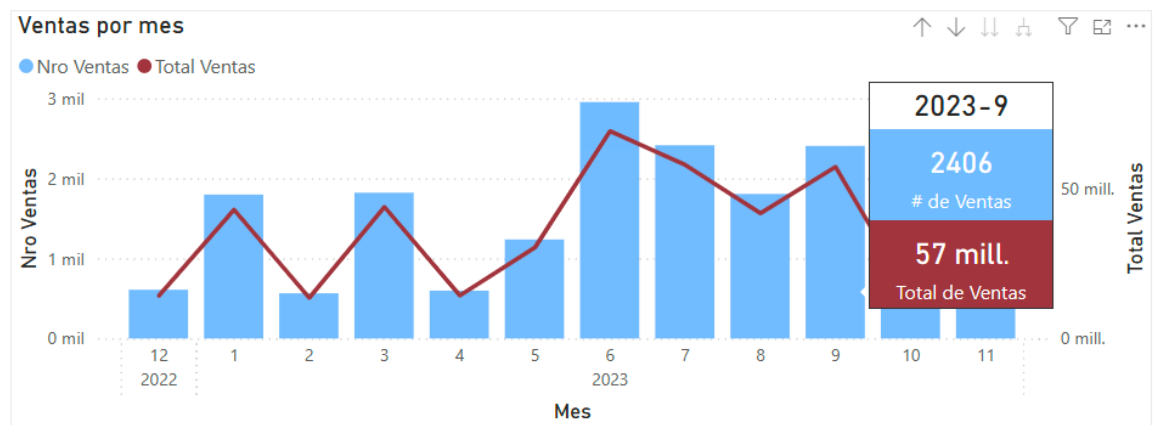
- Si el usuario ubica el mouse sobre la gráfica en una fecha, automáticamente aparecerá una ventana emergente con la información correspondiente a esa fecha:

- En primer lugar, aparecerá la fecha en formato dd/mm/yyyy.
- En color azul, aparecerá la cantidad de Ventas realizadas en la fecha.
- En color rojo, aparecerá el Total de Ventas realizadas en la fecha.

representan la cantidad de ventas mensuales realizadas y la línea representa el valor total de Ventas mensuales.

Si el usuario coloca el mouse sobre la gráfica en una fecha, automáticamente aparecerá una ventana emergente con la información correspondiente a esa fecha:

- En primer lugar, aparecerá la fecha en formato yyyy-mm.
- En color azul, aparecerá la cantidad de Ventas realizadas en el mes.
- En color rojo, aparecerá el Total de Ventas realizadas en el mes.



12. Un gráfico circular que permite analizar el número de ventas realizadas de acuerdo al estado (Entregado o En Camino), para la última fecha de venta, de acuerdo al filtro aplicado por el usuario

Conclusiones

La limpieza y transformación de datos es fundamental para obtener insights precisos y confiables. Al garantizar la consistencia y precisión de la información, hemos podido identificar patrones y tendencias relevantes en los datos de ventas ficticias, lo que a su vez ha permitido tomar decisiones más informadas.

La visualización de datos es una herramienta poderosa para comunicar hallazgos complejos de manera clara y efectiva. El dashboard interactivo desarrollado en Power BI ha permitido explorar los datos de múltiples formas, revelando insights que de otra manera podrían haber pasado desapercibidos. Esto demuestra la importancia de combinar técnicas de limpieza de datos con herramientas de visualización para obtener un análisis completo.

La práctica con datos artificiales es una excelente manera de desarrollar habilidades en limpieza de datos y creación de visualizaciones. Al trabajar con un conjunto de datos generado artificialmente, hemos podido experimentar con diferentes técnicas y herramientas, lo que nos ha permitido fortalecer nuestras competencias en este campo.