

Workshop - 001: Data engineer

Por: Juan Pablo Gómez Veira

Carrera: Ingeniería de Datos e Inteligencia Artificial

Materia: ETL

Introducción

Este documento describe el desarrollo del Workshop 001 , un ejercicio práctico que simula un desafío real de una entrevista de trabajo. El objetivo fue implementar un proceso completo de Extracción, Transformación y Carga (ETL) utilizando datos de candidatos generados aleatoriamente, almacenarlos en una base de datos relacional y finalmente crear visualizaciones específicas en un dashboard interactivo.

Para llevar a cabo el workshop, utilicé las siguientes tecnologías:

- **Python:** Para manipular y transformar los datos.
- **Jupyter Notebook:** Para documentar y ejecutar el código de manera organizada.
- **PostgreSQL:** Como base de datos relacional para almacenar la información.
- **Power BI:** Para diseñar el dashboard final con las visualizaciones solicitadas.

La consigna del workshop pedía migrar los datos desde un archivo CSV (candidates.csv) con 50,000 filas y 10 columnas a una base de datos, aplicar transformaciones —como determinar qué candidatos fueron contratados según una regla específica— y mostrar métricas en cuatro tipos de gráficos: un pie chart, un bar chart horizontal, un bar chart vertical y un multiline chart.

Todo el proceso está disponible en el siguiente repositorio de GitHub:
https://github.com/jpgomezv/ETL_Workshop_001.

A continuación, detallo cada etapa del proyecto, desde la configuración inicial hasta las conclusiones finales.

Configuración del entorno

Antes de empezar con el ETL, preparé un entorno de desarrollo limpio y reproducible. Esto fue clave para asegurarme de que todo funcionara correctamente y de que el proyecto fuera fácil de replicar.

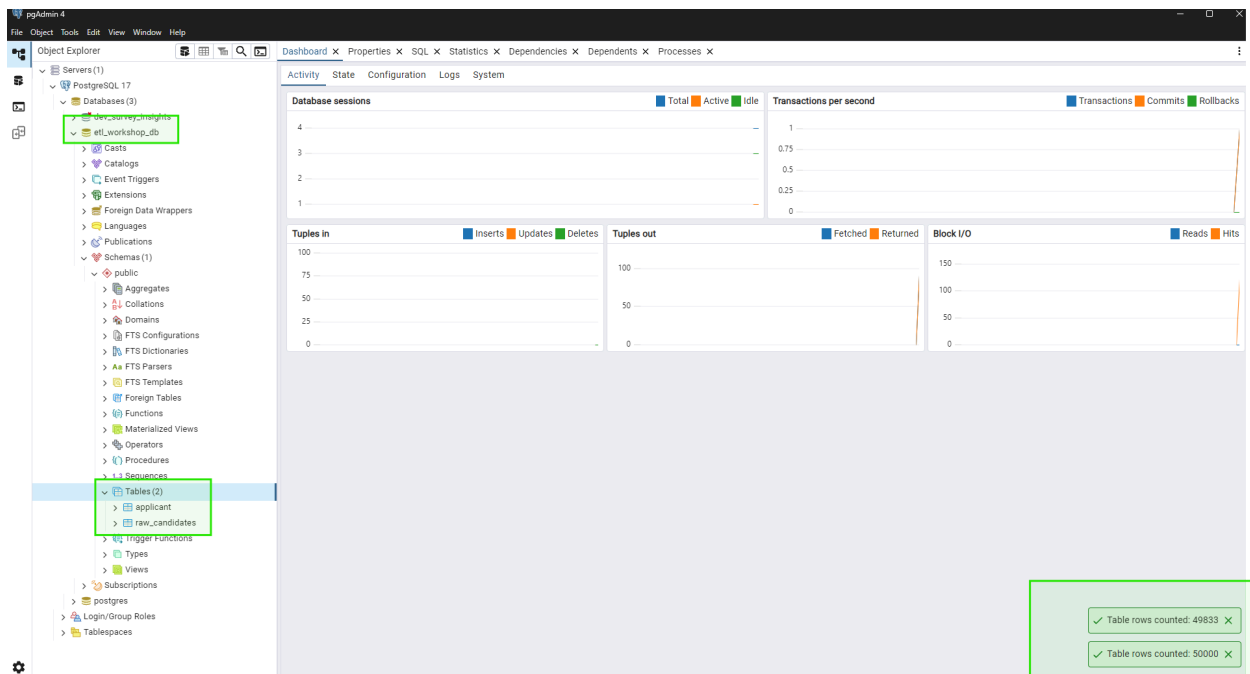
Entorno Virtual y Dependencias

Primero, creé un entorno virtual en Python “venv” para aislar las dependencias del proyecto. Luego, instalé las librerías necesarias usando el archivo requirements.txt, que incluye herramientas como pandas para manipular datos, sqlalchemy para conectar con la base de datos y psycopg2 para trabajar con PostgreSQL.

```
(venv) c:\College\ETL\workshop>pip install -r requirements.txt
Requirement already satisfied: pandas==2.2.3 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 1)) (2.2.3)
Requirement already satisfied: numpy==2.2.3 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 2)) (2.2.3)
Requirement already satisfied: matplotlib==3.10.1 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 3)) (3.10.1)
Requirement already satisfied: seaborn==0.13.2 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 4)) (0.13.2)
Requirement already satisfied: psycopg2-binary==2.9.10 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 5)) (2.9.10)
Requirement already satisfied: SQLAlchemy==2.0.38 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 6)) (2.0.38)
Requirement already satisfied: python-dotenv==1.0.1 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 7)) (1.0.1)
Requirement already satisfied: jupyter==1.1.1 in c:\college\etl\workshop\venv\lib\site-packages (from -r requirements.txt (line 8)) (1.1.1)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\college\etl\workshop\venv\lib\site-packages (from pandas==2.2.3->-r requirements.txt (line 1)) (2.9.0.post0)
Requirement already satisfied: pytz>=2020.1 in c:\college\etl\workshop\venv\lib\site-packages (from pandas==2.2.3->-r requirements.txt (line 1)) (2025.1)
Requirement already satisfied: tzdata>=2022.7 in c:\college\etl\workshop\venv\lib\site-packages (from pandas==2.2.3->-r requirements.txt (line 1)) (2025.1)
Requirement already satisfied: contourpy>=1.0.1 in c:\college\etl\workshop\venv\lib\site-packages (from matplotlib==3.10.1->-r requirements.txt (line 3)) (1.3.1)
Requirement already satisfied: cycler>=0.10 in c:\college\etl\workshop\venv\lib\site-packages (from matplotlib==3.10.1->-r requirements.txt (line 3)) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in c:\college\etl\workshop\venv\lib\site-packages (from matplotlib==3.10.1->-r requirements.txt (line 3)) (4.56.0)
Requirement already satisfied: kiwisolver>=1.3.1 in c:\college\etl\workshop\venv\lib\site-packages (from matplotlib==3.10.1->-r requirements.txt (line 3)) (1.4.8)
Requirement already satisfied: packaging>=20.0 in c:\college\etl\workshop\venv\lib\site-packages (from matplotlib==3.10.1->-r requirements.txt (line 3)) (24.2)
Requirement already satisfied: pillow>=8 in c:\college\etl\workshop\venv\lib\site-packages (from matplotlib==3.10.1->-r requirements.txt (line 3)) (11.1.0)
Requirement already satisfied: pyparsing>=2.3.1 in c:\college\etl\workshop\venv\lib\site-packages (from matplotlib==3.10.1->-r requirements.txt (line 3)) (3.2.1)
```

Configuración de la Base de Datos

Elegí PostgreSQL como base de datos relacional por su robustez y facilidad de integración con Python. Creé una base de datos llamada etl_workshop_db y dos tablas: raw_candidates para los datos crudos y applicant para los datos transformados. Para automatizar este paso, escribí un script en setup.py que configura todo desde cero.



Dividí el proceso ETL entre los notebooks `01_extract_load.ipynb` y `03_transform_data.ipynb`

Extracción

En esta etapa, cargué los datos del archivo `candidates.csv` a la tabla `raw_candidates` en PostgreSQL. Usé `pandas` para leer el CSV y `sqlalchemy` para establecer la conexión con la base de datos y subir los datos. El archivo contenía 50,000 filas con columnas como `first_name`, `email`, `country`, `application_date`, `technology`, entre otras.

CSV File Reading

This section reads the `candidates.csv` file using `pandas`. It includes error handling to ensure the process does not fail if something goes wrong.

```
try:
    df = pd.read_csv('../data/candidates.csv', sep=';')
except FileNotFoundError:
    print("Error: candidates.csv not found")
    raise

# Validate expected columns
expected_columns = ['First Name', 'Last Name', 'Email', 'Application Date', 'Country', 'YOE', 'Seniority', 'Technology', 'Code Challenge Score', 'Technical Interview Score']
if list(df.columns) != expected_columns:
    print("Error: CSV does not have the expected columns.")
    raise ValueError("Column mismatch in candidates.csv")

[28] ✓ 0.0s

if not df.empty:
    print(f"Data loaded successfully. Rows: {len(df)}, Columns: {len(df.columns)}")
else:
    print("The DataFrame is empty.")

[29] ✓ 0.0s
... Data loaded successfully. Rows: 50000, Columns: 10
```

Data Loading to the Database

The data is loaded into the `raw_candidates` table of PostgreSQL replacing the table if it already exists.

```
try:
    df.to_sql('raw_candidates', engine, if_exists='append', index=False)
    print("Data loaded into raw_candidates.")
except Exception as e:
    print(f"Error loading data into PostgreSQL: {e}")

[30] ✓ 1.5s
... Data loaded into raw_candidates.
```

Verification of Data in raw_candidates

The following section displays the first 5 rows of the `raw_candidates` table to confirm that the data was loaded correctly:

```
# Verify that everything is correct
query = "SELECT * FROM raw_candidates LIMIT 5;"
pd.read_sql(query, engine)

[31] ✓ 0.0s
...


|   | id | First Name | Last Name  | Email                     | Application Date | Country | YOE | Seniority | Technology                        | Code Challenge Score | Technical Interview Score |
|---|----|------------|------------|---------------------------|------------------|---------|-----|-----------|-----------------------------------|----------------------|---------------------------|
| 0 | 1  | Bernadette | Langworth  | leonard91@yahoo.com       | 2021-02-26       | Norway  | 2   | Intern    | Data Engineer                     | 3                    | 3                         |
| 1 | 2  | Camryn     | Reynolds   | zelda56@hotmail.com       | 2021-09-09       | Panama  | 10  | Intern    | Data Engineer                     | 2                    | 10                        |
| 2 | 3  | Larue      | Spinka     | okey_schultz41@gmail.com  | 2020-04-14       | Belarus | 4   | Mid-Level | Client Success                    | 10                   | 9                         |
| 3 | 4  | Arch       | Spinka     | elvera_kulas@yahoo.com    | 2020-10-01       | Eritrea | 25  | Trainee   | QA Manual                         | 7                    | 1                         |
| 4 | 5  | Larue      | Altenwerth | minnie.gislason@gmail.com | 2020-05-20       | Myanmar | 13  | Mid-Level | Social Media Community Management | 9                    | 7                         |


```

(Nota: en esa captura en "if_exists='append'", va es **replace**, aclaro por que cuando lo tenía en append se estaban cargando los datos varias veces y esto claramente es inconveniente)

Transformación

Las transformaciones las realicé en el notebook `03_transform_data.ipynb`. Aquí apliqué varias operaciones para limpiar y preparar los datos:

- Creé la columna hired, que indica si un candidato fue contratado. Según la consigna, un candidato es contratado si su code_challenge_score y technical_interview_score son ambos mayores o iguales a 7.
- Eliminé duplicados en la columna email, quedándome con el registro más reciente basado en application_date.
- Normalicé las columnas technology y country para corregir inconsistencias (por ejemplo, mayúsculas o espacios).
- Agrupé las tecnologías en categorías más amplias para simplificar el análisis.

```

Data Transformations

In this section, we transform the data to create the final applicant table. This includes creating a boolean hired column based on the criteria identified in the EDA (scores >= 7 in both code_challenge_score and technical_interview_score), removing duplicate emails (if applicable), normalizing data, and grouping technologies into broader categories to facilitate analysis.

# Create the 'hired' column based on EDA criteria
df['hired'] = ((df['code_challenge_score'] >= 7) & (df['technical_interview_score'] >= 7)).astype(bool)

# Remove duplicate emails, keeping the most recent record (based on application_date)
df = df.sort_values('application_date', ascending=False).drop_duplicates(subset='email', keep='first')

# Normalize 'technology' and 'country' columns
df['technology'] = df['technology'].str.strip().str.title()
df['country'] = df['country'].str.strip().str.title()

# Dictionary to map technologies into broader categories
technology_mapping = {
    'Data Engineer': 'Data & Analytics',
    'Client Success': 'Customer Support',
    'Qa Manual': 'Quality Assurance (QA)',
    'Social Media Community Management': 'Marketing & Community',
    'Adobe Experience Manager': 'Marketing & Community',
    'Sales': 'Sales',
    'Salesforce': 'Customer Relationship Mgmt',
    'System Administration': 'DevOps & Infrastructure',
    'Security': 'Security',
    'Game Development': 'Development - Other',
    'Development - Cms Frontend': 'Development - Frontend',
    'Security Compliance': 'Security',
    'Development - Backend': 'Development - Backend',
    'Design': 'Design',
    'Business Analytics / Project Management': 'Business & Project Management',
    'Development - Frontend': 'Development - Frontend',

```

```
# Apply the mapping to the 'technology' column
df['technology'] = df['technology'].map(technology_mapping)

# Verify the results
print("Number of candidates after removing duplicate emails:", len(df))
print("Percentage of hired candidates:", (df['hired'].mean() * 100).round(2), "%")
print("\nUnique technology categories after grouping:", df['technology'].unique())
df.head()
```

✓ 0.1s

Number of candidates after removing duplicate emails: 49833
Percentage of hired candidates: 13.41 %

Unique technology categories after grouping: ['Customer Relationship Mgmt' 'Design' 'Documentation & Writing'
'Data & Analytics' 'DevOps & Infrastructure' 'Marketing & Community'
'Business & Project Management' 'Quality Assurance (QA)'
'Development - Other' 'Customer Support' 'Sales' 'Development - Frontend'
'Security' 'Development - FullStack' 'Development - Backend'
'Integration & Middleware']

Carga

Una vez transformados, cargué los datos finales en la tabla applicant de PostgreSQL. Esta tabla contiene la versión limpia y lista para análisis, con la columna hired y los ajustes necesarios.

```
# Validate the 'hired' column for null values
if df['hired'].isnull().sum() > 0:
    print("Warning: There are null values in the 'hired' column.")

try:
    df.to_sql('applicant', engine, if_exists='append', index=False)
    print("Transformed data successfully loaded into the applicant table.")
except Exception as e:
    print(f"Error loading data into PostgreSQL: {e}")
```

✓ 19s

Transformed data successfully loaded into the applicant table.

Verification of the Final Table

We confirm that the transformed data was correctly loaded into the `applicant` table in PostgreSQL.

```
# Verify the data in the applicant table
query = "SELECT * FROM applicant LIMIT 5;"
final_df = pd.read_sql(query, engine)
print("First 5 records of the applicant table:")
final_df
```

✓ 0.0s

First 5 records of the applicant table:

	id	first_name	last_name	email	application_date	country	years_of_experience	seniority	technology	code_challenge_score	technical_interview_score	hired
0	39931	Cleveland	Ledner	kiana.lang86@gmail.com	2022-07-04	Iraq	12	Trainee	Data & Analytics	8	6	False
1	34859	Myles	Towne	kris.green@hotmail.com	2022-07-04	Kyrgyz Republic	18	Senior	Business & Project Management	7	9	True
2	23099	Philip	Hessel	kallie.hilpert48@gmail.com	2022-07-04	Uzbekistan	28	Intern	Customer Support	4	5	False
3	45937	Maritza	Beer	theodora.crimes28@gmail.com	2022-07-04	Algeria	10	Architect	Data & Analytics	6	8	False

Análisis Exploratorio de Datos (EDA)

Antes de pasar a las visualizaciones finales, hice un análisis exploratorio en el notebook `02_explore_data.ipynb` para entender mejor los datos y detectar patrones o problemas.

4.1. Hallazgos Clave

- No encontré valores nulos en el dataset, lo cual facilitó el proceso.
- Había duplicados en email, pero los manejé en la etapa de transformación.
- Las variables categóricas como seniority y technology mostraron distribuciones variadas; por ejemplo, algunas tecnologías eran más comunes que otras.
- Las puntuaciones code_challenge_score y technical_interview_score tenían una distribución uniforme, típico de datos generados aleatoriamente. No hay una relación explicada entre ambos puntajes
- Al aplicar la regla de contratación, descubrí que aproximadamente el 13.4% de los candidatos fueron contratados.

4.2. Visualizaciones del EDA

Generé varias gráficas para explorar los datos, como por ejemplo:

- Un histograma de la distribución de seniority.
- Un conteo de las tecnologías más frecuentes.
- Un heatmap para ver la relación entre code_challenge_score y technical_interview_score.

```
Hired Candidate Analysis

This section temporarily explores candidates meeting the criteria for being considered 'hired' (scores of 7 or higher in both code challenge and technical interview). This is an exploratory analysis and does not involve permanent data transformations, which will be handled in the transformation notebook.

# Filter 'hired' candidates (inspired by the image and temp_eda_report.ipynb)
df_hired = df.loc[(df['code_challenge_score'] >= 7) & (df['technical_interview_score'] >= 7)]

# Display a summary of hired candidates
hired_count = len(df_hired)
total_count = len(df)
print(f"Number of hired candidates: {hired_count} ({(hired_count / total_count * 100):.2f}% of total)")

Number of hired candidates: 13396 (13.48% of total)
```

```

# Check for null values
null_counts = df.isnull().sum()
print("Null values by column:\n", null_counts)

# Check for duplicate rows
duplicates = df.duplicated().sum()
print(f"Number of duplicate rows: {duplicates}")

# Check for duplicate emails
email_duplicates = df['email'].duplicated().sum()
print(f"Number of duplicate emails (excluding the first): {email_duplicates}")

# Review duplicate records
email_duplicates = df.duplicated(subset='email', keep=False).sum()
print("Records sharing the same email:", email_duplicates)
# Optionally display some duplicate email records
if email_duplicates > 0:
    duplicated_emails_df = df[df.duplicated(subset='email', keep=False)].sort_values('email')
    display(duplicated_emails_df.head(10))

# Note on duplicate emails (given the fictional context)
print("""
Note: Duplicate emails might suggest re-applications in a real context, but since the data is randomly generated, \nthese repetitions are likely artifacts of the generation process and may not i
""")

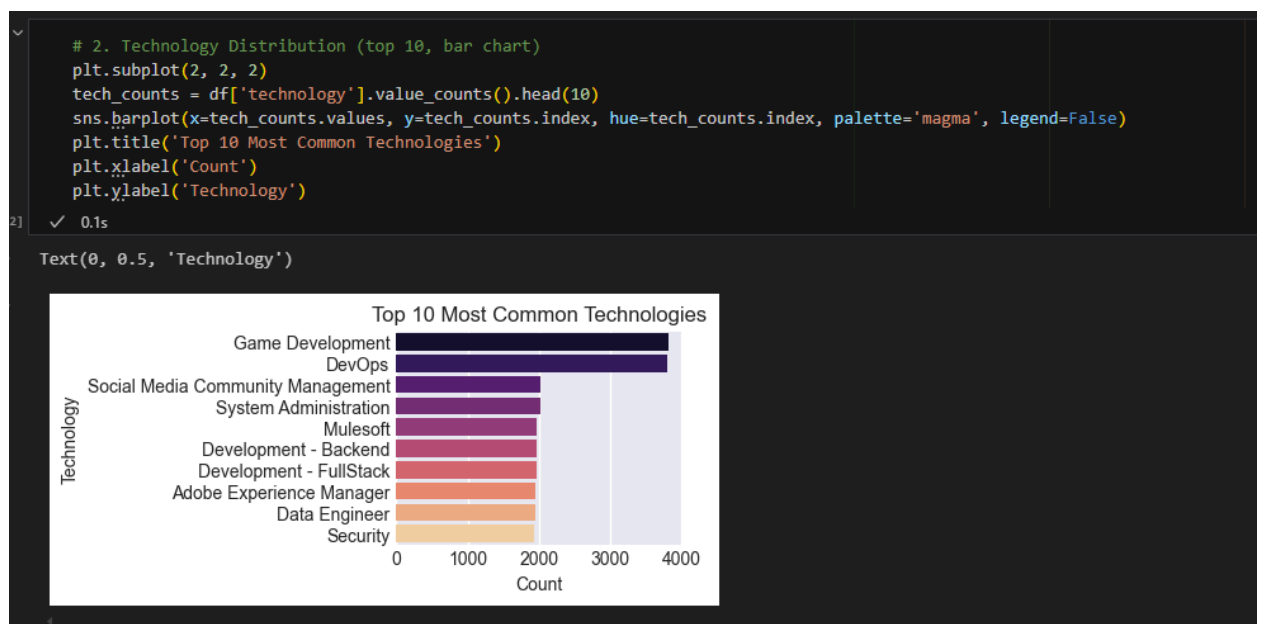
```

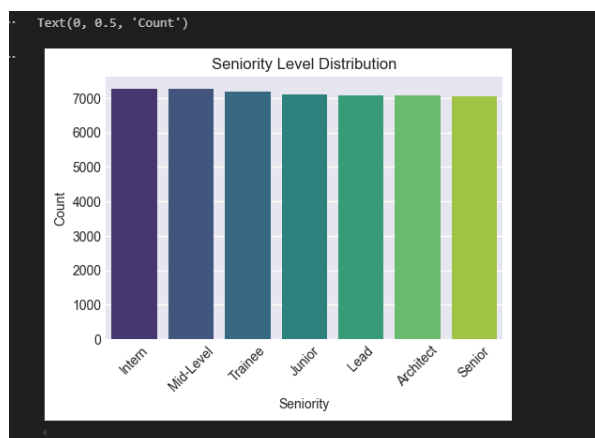
✓ 0.0s

```

Null values by column:
first_name      0
last_name       0
email           0
application_date 0
country         0
years_of_experience 0
seniority        0
technology       0
code_challenge_score 0
technical_interview_score 0
dtype: int64
Number of duplicate rows: 0
Number of duplicate emails (excluding the first): 167
Records sharing the same email: 332

```



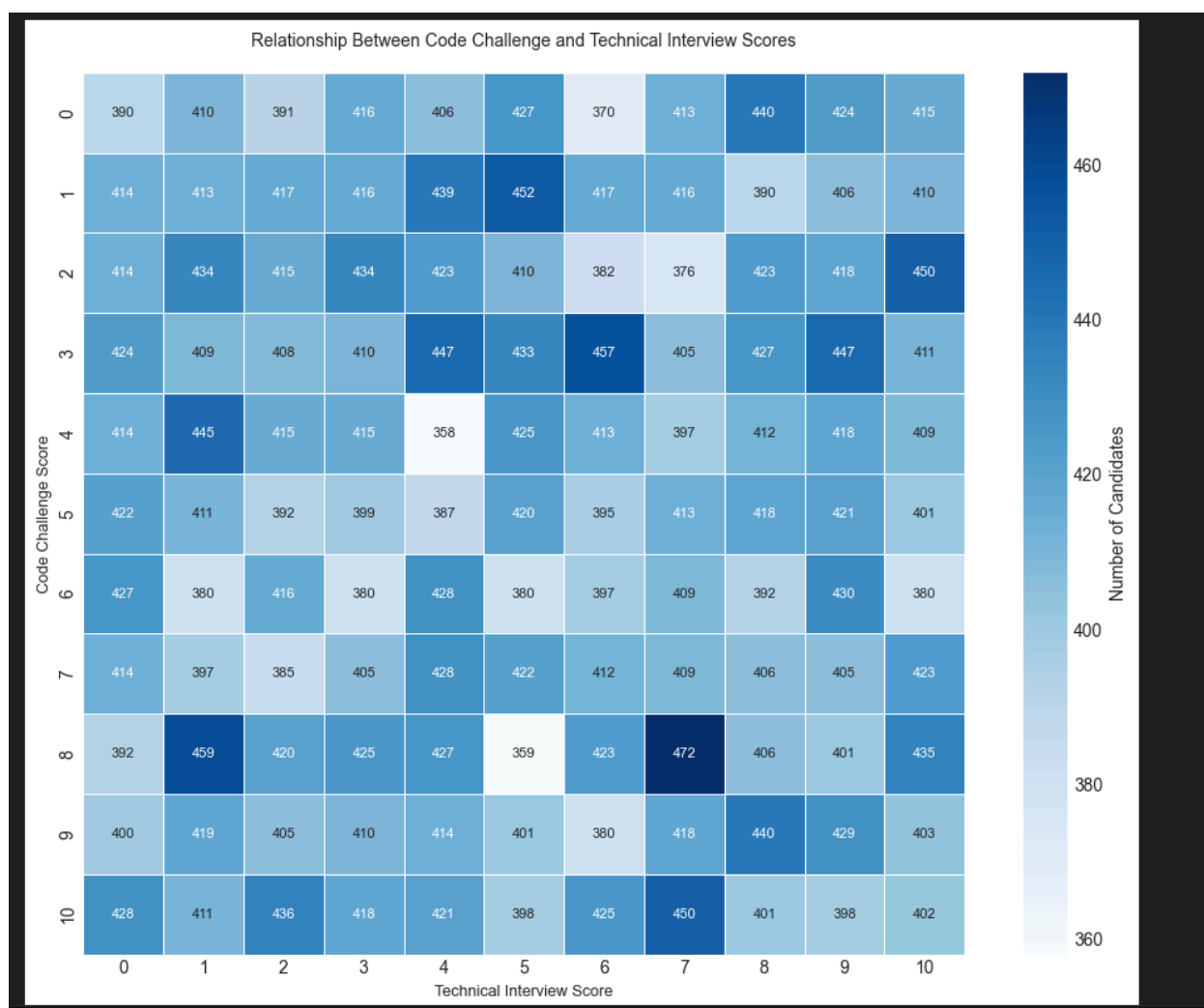


Distribution of seniority:

```

seniority
Intern      7255
Mid-Level   7253
Trainee     7183
Junior      7100
Architect   7079
Lead        7071
Senior      7059
Name: count, dtype: int64

```



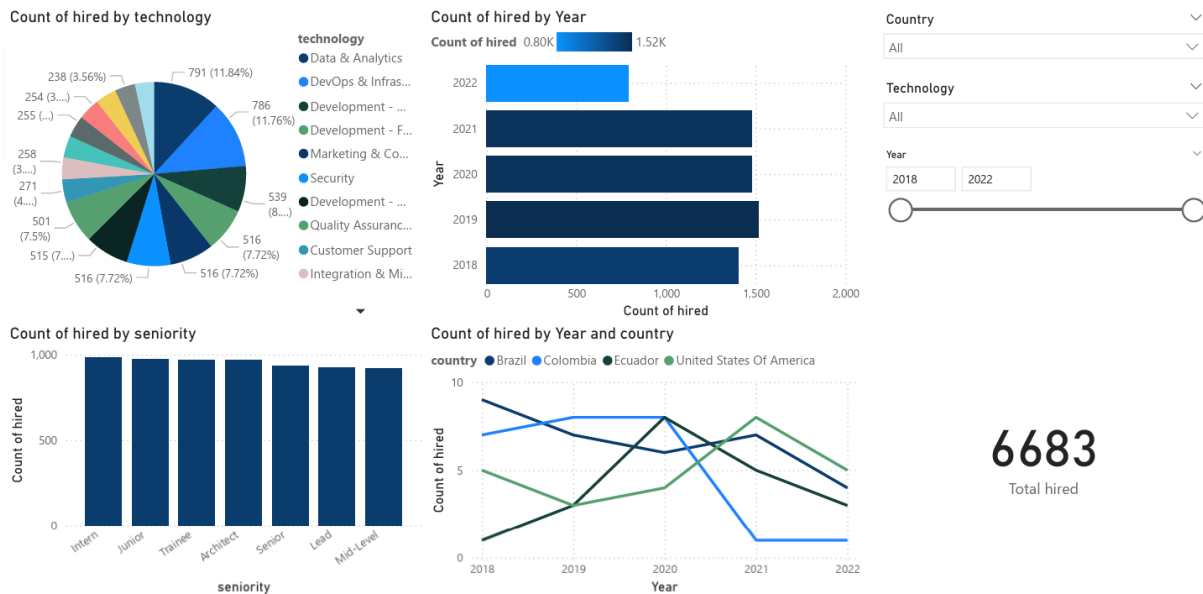
(Dentro del 02_explore_data.ipynb se encuentran más graficas, pero con estas ya se entienden los hallazgos encontrados)

Visualizaciones Finales

Para las visualizaciones finales, usé Power BI conectado directamente a la tabla applicant en PostgreSQL. Diseñé un dashboard con los cuatro gráficos solicitados en la consigna:

1. **Contrataciones por Tecnología (Pie Chart):** Muestra cómo se distribuyen los candidatos contratados entre las categorías de tecnología.
2. **Contrataciones por Año (Horizontal Bar Chart):** Presenta el número de contrataciones por año entre 2018 y 2022.
3. **Contrataciones por Seniority (Bar Chart):** Indica cuántos candidatos contratados hay por nivel de experiencia.
4. **Contrataciones por País a lo Largo de los Años (Multiline Chart):** Muestra las tendencias de contrataciones en USA, Brasil, Colombia y Ecuador a lo largo del tiempo.

Añadí filtros interactivos para Country, Technology y Year, lo que permite explorar los datos de forma dinámica.



Conclusiones

Este workshop fue una gran oportunidad para poner a prueba practicas que se deben tener en proyectos de ETL. Desde configurar el entorno hasta transformar datos y crear visualizaciones, cada paso me enseñó algo nuevo. Algunos puntos clave que rescato:

- Automatizar la configuración del entorno con scripts como setup.py ahorra tiempo y asegura consistencia.
- El EDA es esencial para conocer los datos antes de trabajar con ellos.
- Aplicar reglas de negocio, como la de contratación, requiere pensar bien cómo traducirlas a código.
- Power BI es una herramienta poderosa para presentar resultados de manera clara y profesional (en lo personal me gustó la simplicidad y solides de uso).

Si tuviera que mejorar el proyecto, intentaría:

- Agregar pruebas para validar las transformaciones.
- Optimizar el código para manejar datasets más grandes.
- Explorar análisis adicionales, como predecir tendencias de contratación.

