**Visualización para grandes volúmenes de datos**

**Mario Andrés Hernández Moreno**

**Actividad 2: Proyecto Integrador de Sistemas de Modelado, Procesamiento y Gestión**

**Mag. Mario Alejandro Bravo Ortiz**

**Universidad Autónoma de Manizales, Manizales**

**Especialización en Inteligencia Artificial**

**2024**

**Contenido**

# 1. Reseña escrita

El modelado y la gestión eficiente de grandes volúmenes de datos se han vuelto cruciales en la era digital actual, pues a medida que las organizaciones generan y recolecta grandes cantidades de información, la capacidad para tratar, procesar y extraer valor de estos datos se ha convertido en una necesidad y en una ventaja competitiva clave para el mundo empresarial y el desarrollo de diversos sectores.

El modelado de datos permite estructurar y organizar la información de manera que sea fácilmente accesible y analizable, lo que facilita la identificación de patrones, tendencias, correlaciones y nuevos puntos de vista que pueden impulsar la toma de decisiones con información valiosa y también la innovación. Por otro lado, una gestión eficiente garantiza que los datos se almacenen, procesen y transmitan de manera segura y rentable preservando su integridad y seguridad, esto se traduce en una ventaja competitiva clave, permitiendo a las empresas mejorar sus operaciones, personalizar productos y servicios, y anticipar tendencias futuras.

Sin embargo, el manejo de Big Data también presenta desafíos importantes que se deben superar, como lo son la integración de datos heterogéneos, el procesamiento en tiempo real, la necesidad de infraestructura adecuada, la protección de la privacidad y la garantía de la calidad de los datos. Abordar estas dificultades es esencial y fundamental para aprovechar al máximo el potencial de los grandes volúmenes de datos en nuestra sociedad y aplicarlos de una manera adecuada en el día a día de las organizaciones.

En resumen, el modelado y la gestión eficiente de grandes volúmenes de datos son fundamentales para convertir datos brutos en información valiosa, facilitando la toma de decisiones informadas, impulsando la innovación y promoviendo el progreso en diversos sectores. Su importancia radica en su capacidad para transformar la manera en que las empresas operan, investigan y evolucionan en un mundo cada vez más interesado en la información y con una dependencia creciente hacia el Big Data.

## 2. Airbnb Listing Data 2023



Configuración del ambiente de Google Colaboratory con Pyspark y Hadoop

```
# Download Java
!apt-get install openjdk-8-jdk-headless -qq > /dev/null
# Next, we will install Apache Spark 3.0.1 with Hadoop 2.7 from here.
!wget https://archive.apache.org/dist/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz
# Now, we just need to unzip that folder.
!tar xf spark-3.5.1-bin-hadoop3.tgz

# Setting JVM and Spark path variables
import os
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
os.environ["SPARK_HOME"] = "/content/spark-3.5.1-bin-hadoop3"

# Installing required packages
!pip install pyspark==3.5.1
!pip install findspark
```

Importamos librerías

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import os
import plotly.express as px
import plotly.io as pio

import findspark
findspark.init()

from pyspark import SparkContext
from pyspark.sql import Window
from pyspark.sql import SparkSession
from pyspark.sql.types import DoubleType, IntegerType, DateType
from pyspark.sql import functions as fct
from pyspark.ml  import Pipeline
from pyspark.sql import SQLContext
from pyspark.sql.functions import mean,col,split,regexp_extract,when,lit,sum,desc,to_date,expr,round,avg,max,count,regexp_replace, trim
from pyspark.ml.feature import StringIndexer, VectorAssembler
from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.ml.feature import QuantileDiscretizer
```

Creamos la sesión de Spark

```
[3]  SpSession = SparkSession \
              .builder \
              .appName("Airbnb Spark") \
              .getOrCreate()


[4]  SpContext = SpSession.sparkContext
```

Importamos los datos y los visualizamos

```python
df1 = SpSession.read.csv('/content/airnb.csv', header=True, sep=",")
df1.show(5)

df2 = SpSession.read.csv('/content/airnb_desert.csv', header=True, sep=",")
df2.show(5)

df3 = SpSession.read.csv('/content/airnb_luxe.csv', header=True, sep=",")
df3.show(5)
```

| Title | Detail | Date | Price(in dollar) | Offer price(in dollar) | Review and rating | Number of bed |
|---|---|---|---|---|---|---|
| Chalet in Skykomi... | Sky Haus - A-Fram... | Jun 11 - 16 | 306.00 | 229.00 | 4.85 (531) | 4 beds |
| Cabin in Hancock,... | The Catskill A-Fr... | Jun 6 - 11 | 485.00 | 170.00 | 4.77 (146) | 4 beds |
| Cabin in West Far... | The Triangle: A-F... | Jul 9 - 14 | 119.00 | 522.00 | 4.91 (515) | 4 beds |
| Home in Blue Ridg... | *Summer Sizzle* 5... | Jun 11 - 16 | 192.00 | 348.00 | 4.94 (88) | 5 beds |
| Treehouse in Gran... | Luxury Treehouse ... | Jun 4 - 9 | 232.00 | 196.00 | 4.99 (222) | 1 queen bed |

only showing top 5 rows

| Desert name | Date | Price(In dollar) | Details | Rating |
|---|---|---|---|---|
| Mhamid, Morocco | May 1 ◆ 29 | 479.00 | Near Sahara Desert | 4.79 |
| Aqaba City, Jordan | May 1 ◆ 29 | 2,168.00 | Near Arabian Desert | 4.92 |
| Tamesluht, Morocco | May 1 ◆ 29 | 17,752.00 | 9,404 kilometers ... | 4.95 |
| Al Bairat, Egypt | May 1 ◆ 29 | 1,982.00 | Near Sahara Desert | 4.88 |
| Tamesluht, Morocco | May 1 ◆ 29 | 17,752.00 | 9,404 kilometers ... | 4.87 |

only showing top 5 rows

| Luxe name | Date | Price(In dollar) | Distance |
|---|---|---|---|
| Koh Samui, Thailand | May 1 ◆ 29 | 89,600.00 | 1,880 kilometers ... |
| Koh Samui, Thailand | May 1 ◆ 29 | 78,459.00 | 1,880 kilometers ... |
| Koh Samui, Thailand | May 1 ◆ 29 | 53,200.00 | 1,881 kilometers ... |
| Koh Samui, Thailand | May 1 ◆ 29 | 35,000.00 | 1,880 kilometers ... |
| Nathon, Thailand | May 1 ◆ 29 | 19,656.00 | 1,872 kilometers ... |

only showing top 5 rows

Exploración de esquemas

```
root
 |-- Title: string (nullable = true)
 |-- Detail: string (nullable = true)
 |-- Date: string (nullable = true)
 |-- Price(in dollar): string (nullable = true)
 |-- Offer price(in dollar): string (nullable = true)
 |-- Review and rating: string (nullable = true)
 |-- Number of bed: string (nullable = true)

root
 |-- Desert name: string (nullable = true)
 |-- Date: string (nullable = true)
 |-- Price(In dollar): string (nullable = true)
 |-- Details: string (nullable = true)
 |-- Rating: string (nullable = true)

root
 |-- Luxe name: string (nullable = true)
 |-- Date: string (nullable = true)
 |-- Price(In dollar): string (nullable = true)
 |-- Distance: string (nullable = true)
```

Características estadísticas básicas

| summary | Title | Detail | Date | Price(in dollar) | Offer price(in dollar) | Review and rating | Number of bed |
|---|---|---|---|---|---|---|---|
| count | 957 | 957 | 953 | 949 | 172 | 951 | 949 |
| mean | NULL | NULL | 99.0 | 170.51797040169134 | 149.90419161676647 | NULL | NULL |
| stddev | NULL | NULL | NULL | 88.1702897806285 | 135.57467024158484 | NULL | NULL |
| min | A Resort by the Sea" | """ La cabane du ... | rural tourism. Y... | 1,463.00 | 1,089.00 | 1 sofa bed | 1 bed |
| max | Yurt in Rising Fawn | 📍Boka Bay view a... | Sep 9 - 14 | Sep 16 - 21 | 99.00 | None | 9 beds |

| summary | Desert name | Date | Price(In dollar) | Details | Rating |
|---|---|---|---|---|---|
| count | 280 | 280 | 280 | 280 | 280 |
| mean | NULL | NULL | 554.92 | NULL | 4.721904761904762 |
| stddev | NULL | NULL | 251.24331478828736 | NULL | 0.2559845438376616 |
| min | ???? ???? ??, Jordan | May 1 ◆ 29 | 1,041.00 | 10,017 kilometers... | 3.5 |
| max | ◆rzola, Spain | May 5 ◆ Jun 2 | 998.00 | Near Thar Desert | None |

| summary | Luxe name | Date | Price(In dollar) | Distance |
|---|---|---|---|---|
| count | 280 | 280 | 280 | 280 |
| mean | NULL | NULL | NULL | NULL |
| stddev | NULL | NULL | NULL | NULL |
| min | A. Thalang, Thailand | May 1 ◆ 29 | 100,047.00 | 1,860 kilometers ... |
| max | Ysterni, Greece | May 6 ◆ Jun 3 | 96,793.00 | 9,954 kilometers ... |

Preparación de los datos, consultas y transformaciones:

Identificación de valores nulos

```
Valores nulos en 'Title': 0
Valores nulos en 'Detail': 0
Valores nulos en 'Date': 4
Valores nulos en 'Price(in dollar)': 8
Valores nulos en 'Offer price(in dollar)': 785
Valores nulos en 'Review and rating': 6
Valores nulos en 'Number of bed': 8


Valores nulos en 'Desert name': 0
Valores nulos en 'Date': 0
Valores nulos en 'Price(In dollar)': 0
Valores nulos en 'Details': 0
Valores nulos en 'Rating': 0


Valores nulos en 'Luxe name': 0
Valores nulos en 'Date': 0
Valores nulos en 'Price(In dollar)': 0
Valores nulos en 'Distance': 0
```

Eliminación de la columna 'OfferPrice' de df1 debido a que no es relevante para nuestro análisis

| Title | Detail | Date | Price(in dollar) | Review and rating | Number of bed |
|---|---|---|---|---|---|
| Chalet in Skykomi... | Sky Haus - A-Fram... | Jun 11 - 16 | 306.00 | 4.85 (531) | 4 beds |
| Cabin in Hancock,... | The Catskill A-Fr... | Jun 6 - 11 | 485.00 | 4.77 (146) | 4 beds |
| Cabin in West Far... | The Triangle: A-F... | Jul 9 - 14 | 119.00 | 4.91 (515) | 4 beds |

only showing top 3 rows

Renombramos columnas y etiquetamos

```
+--------------------+--------------------+-----------+------+-----------------+-------------+--------+
|                Name|              Detail|       Date| Price|Review and rating|Number of bed|Category|
+--------------------+--------------------+-----------+------+-----------------+-------------+--------+
|Chalet in Skykomi...|Sky Haus - A-Fram...|Jun 11 - 16|306.00|         4.85 (531)|       4 beds|Standard|
|Cabin in Hancock,...|The Catskill A-Fr...| Jun 6 - 11|485.00|         4.77 (146)|       4 beds|Standard|
|Cabin in West Far...|The Triangle: A-F...| Jul 9 - 14|119.00|         4.91 (515)|       4 beds|Standard|
+--------------------+--------------------+-----------+------+-----------------+-------------+--------+
only showing top 3 rows


+-----------------+---------+---------+--------------------+------+--------+
|             Name|     Date|    Price|             Details|Rating|Category|
+-----------------+---------+---------+--------------------+------+--------+
|   Mhamid, Morocco|May 1 � 29|   479.00|  Near Sahara Desert|  4.79|  Desert|
|Aqaba City, Jordan|May 1 � 29| 2,168.00| Near Arabian Desert|  4.92|  Desert|
|Tamesluht, Morocco|May 1 � 29|17,752.00|9,404 kilometers ...|  4.95|  Desert|
+-----------------+---------+---------+--------------------+------+--------+
only showing top 3 rows


+-----------------+---------+---------+--------------------+--------+
|             Name|     Date|    Price|            Distance|Category|
+-----------------+---------+---------+--------------------+--------+
|Koh Samui, Thailand|May 1 � 29|89,600.00|1,880 kilometers ...|    Luxe|
|Koh Samui, Thailand|May 1 � 29|78,459.00|1,880 kilometers ...|    Luxe|
|Koh Samui, Thailand|May 1 � 29|53,200.00|1,881 kilometers ...|    Luxe|
+-----------------+---------+---------+--------------------+--------+
```

Extracción del mes de Date, filtrado y eliminación de valores atípicos

```
+--------------------+--------------------+-----------+------+-----------------+-------------+--------+-----+
|                Name|              Detail|       Date| Price|Review and rating|Number of bed|Category|Month|
+--------------------+--------------------+-----------+------+-----------------+-------------+--------+-----+
|Chalet in Skykomi...|Sky Haus - A-Fram...|Jun 11 - 16|306.00|         4.85 (531)|       4 beds|Standard|  Jun|
|Cabin in Hancock,...|The Catskill A-Fr...| Jun 6 - 11|485.00|         4.77 (146)|       4 beds|Standard|  Jun|
|Cabin in West Far...|The Triangle: A-F...| Jul 9 - 14|119.00|         4.91 (515)|       4 beds|Standard|  Jul|
|Home in Blue Ridg...|*Summer Sizzle* 5...|Jun 11 - 16|192.00|          4.94 (88)|       5 beds|Standard|  Jun|
|Treehouse in Gran...|Luxury Treehouse ...|  Jun 4 - 9|232.00|         4.99 (222)| 1 queen bed|Standard|  Jun|
+--------------------+--------------------+-----------+------+-----------------+-------------+--------+-----+
only showing top 5 rows


+-----------------+---------+---------+--------------------+------+--------+-----+
|             Name|     Date|    Price|             Details|Rating|Category|Month|
+-----------------+---------+---------+--------------------+------+--------+-----+
|   Mhamid, Morocco|May 1 � 29|   479.00|  Near Sahara Desert|  4.79|  Desert|  May|
|Aqaba City, Jordan|May 1 � 29| 2,168.00| Near Arabian Desert|  4.92|  Desert|  May|
|Tamesluht, Morocco|May 1 � 29|17,752.00|9,404 kilometers ...|  4.95|  Desert|  May|
|   Al Bairat, Egypt|May 1 � 29| 1,982.00|  Near Sahara Desert|  4.88|  Desert|  May|
|Tamesluht, Morocco|May 1 � 29|17,752.00|9,404 kilometers ...|  4.87|  Desert|  May|
+-----------------+---------+---------+--------------------+------+--------+-----+
only showing top 5 rows


+-----------------+---------+---------+--------------------+--------+-----+
|             Name|     Date|    Price|            Distance|Category|Month|
+-----------------+---------+---------+--------------------+--------+-----+
|Koh Samui, Thailand|May 1 � 29|89,600.00|1,880 kilometers ...|    Luxe|  May|
|Koh Samui, Thailand|May 1 � 29|78,459.00|1,880 kilometers ...|    Luxe|  May|
|Koh Samui, Thailand|May 1 � 29|53,200.00|1,881 kilometers ...|    Luxe|  May|
|Koh Samui, Thailand|May 1 � 29|35,000.00|1,880 kilometers ...|    Luxe|  May|
|   Nathon, Thailand|May 1 � 29|19,656.00|1,872 kilometers ...|    Luxe|  May|
+-----------------+---------+---------+--------------------+--------+-----+
only showing top 5 rows
```

Casteo de Price como un valor numérico

```
+------------------+-----------------+-----------+-----+-----------------+-------------+--------+-----+
|              Name|           Detail|       Date|Price|Review and rating|Number of bed|Category|Month|
+------------------+-----------------+-----------+-----+-----------------+-------------+--------+-----+
|Chalet in Skykomi...|Sky Haus - A-Fram...|Jun 11 - 16|306.0|        4.85 (531)|       4 beds|Standard|  Jun|
|Cabin in Hancock,...|The Catskill A-Fr...| Jun 6 - 11|485.0|        4.77 (146)|       4 beds|Standard|  Jun|
|Cabin in West Far...|The Triangle: A-F...| Jul 9 - 14|119.0|        4.91 (515)|       4 beds|Standard|  Jul|
|Home in Blue Ridg...|*Summer Sizzle* 5...|Jun 11 - 16|192.0|         4.94 (88)|       5 beds|Standard|  Jun|
|Treehouse in Gran...|Luxury Treehouse ...|  Jun 4 - 9|232.0|        4.99 (222)| 1 queen bed|Standard|  Jun|
+------------------+-----------------+-----------+-----+-----------------+-------------+--------+-----+
only showing top 5 rows

+------------------+--------+-------+-------------------+------+--------+-----+
|              Name|    Date|  Price|            Details|Rating|Category|Month|
+------------------+--------+-------+-------------------+------+--------+-----+
|   Mhamid, Morocco|May 1 � 29|  479.0|  Near Sahara Desert|  4.79|  Desert|  May|
|Aqaba City, Jordan|May 1 � 29| 2168.0| Near Arabian Desert|  4.92|  Desert|  May|
|Tamesluht, Morocco|May 1 � 29|17752.0|9,404 kilometers ...|  4.95|  Desert|  May|
|  Al Bairat, Egypt|May 1 � 29| 1982.0|  Near Sahara Desert|  4.88|  Desert|  May|
|Tamesluht, Morocco|May 1 � 29|17752.0|9,404 kilometers ...|  4.87|  Desert|  May|
+------------------+--------+-------+-------------------+------+--------+-----+
only showing top 5 rows

+------------------+--------+-------+-------------------+--------+-----+
|              Name|    Date|  Price|           Distance|Category|Month|
+------------------+--------+-------+-------------------+--------+-----+
|Koh Samui, Thailand|May 1 � 29|89600.0|1,880 kilometers ...|    Luxe|  May|
|Koh Samui, Thailand|May 1 � 29|78459.0|1,880 kilometers ...|    Luxe|  May|
|Koh Samui, Thailand|May 1 � 29|53200.0|1,881 kilometers ...|    Luxe|  May|
|Koh Samui, Thailand|May 1 � 29|35000.0|1,880 kilometers ...|    Luxe|  May|
|   Nathon, Thailand|May 1 � 29|19656.0|1,872 kilometers ...|    Luxe|  May|
+------------------+--------+-------+-------------------+--------+-----+
only showing top 5 rows
```

Extracción de los países para los df

```
+------------------+
|           Country|
+------------------+
|           Tunisia|
| Dominican Republic|
|          Colombia|
|           Ireland|
|              Cuba|
|            Taiwan|
|       El Salvador|
|            Canada|
|           Germany|
|           Curaçao|
|                US|
|                UK|
|             Spain|
|           Vietnam|
|           Czechia|
|       Philippines|
|           Jamaica|
|       Puerto Rico|
|          Malaysia|
|       Switzerland|
+------------------+
only showing top 20 rows
```

```
+------------+
|     Country|
+------------+
|     Tunisia|
|   Australia|
|      Israel|
|       Aswan|
|       Spain|
|       Egypt|
|        Oman|
|     Morocco|
|      Jordan|
|     Namibia|
|     Morocco|
|       India|
| Saudi Arabia|
|      Greece|
| South Africa|
|         UAE|
+------------+
```

```
+------------+
|     Country|
+------------+
|     Croatia|
|   Australia|
|    Maldives|
|   Sri Lanka|
|   Mauritius|
|      Turkey|
|  Montenegro|
|   Koh Samui|
|       Italy|
|      Sweden|
|     Austria|
|      Norway|
|   Indonesia|
|    Thailand|
|      Greece|
|     Finland|
| South Africa|
+------------+
```

Unión de los df con algunas columnas

```
+------------------+-----+--------+-------+
|              Name|Price|Category|Country|
+------------------+-----+--------+-------+
|Chalet in Skykomi...|306.0|Standard|     US|
|Cabin in Hancock,...|485.0|Standard|     US|
|Cabin in West Far...|119.0|Standard|     US|
|Home in Blue Ridg...|192.0|Standard|     US|
|Treehouse in Gran...|232.0|Standard|     US|
+------------------+-----+--------+-------+
only showing top 5 rows
```

Filtro de propiedades con precio mayor a 1000 USD

```
+----------------------------------+-------+--------+--------+
|Name                              |Price  |Category|Country |
+----------------------------------+-------+--------+--------+
|Home in Moss Beach, California, US|1463.0 |Standard| US     |
|Aqaba City, Jordan                |2168.0 |Desert  | Jordan |
|Tamesluht, Morocco                |17752.0|Desert  | Morocco|
|Al Bairat, Egypt                  |1982.0 |Desert  | Egypt  |
|Tamesluht, Morocco                |17752.0|Desert  | Morocco|
|Mitzpe Ramon, Israel              |2502.0 |Desert  | Israel |
|Tinghir, Morocco                  |1346.0 |Desert  | Morocco|
|Mitzpe Ramon, Israel              |2945.0 |Desert  | Israel |
|Merzouga, Morocco                 |2937.0 |Desert  | Morocco|
|Essaouira, Morocco                |5386.0 |Desert  | Morocco|
|Tzofar, Israel                    |1821.0 |Desert  | Israel |
|Ait Bihi, Morocco                 |1096.0 |Desert  | Morocco|
|Al Bairat, Egypt                  |1670.0 |Desert  | Egypt  |
|�rzola, Spain                     |2396.0 |Desert  | Spain  |
|Aglou, Morocco                    |2792.0 |Desert  | Morocco|
|Marrakech, Morocco                |3561.0 |Desert  | Morocco|
|Tamesluht, Morocco                |12800.0|Desert  | Morocco|
|Arad, Israel                      |3737.0 |Desert  | Israel |
|Mhamid, Morocco                   |1145.0 |Desert  | Morocco|
|???????, Morocco                  |1149.0 |Desert  | Morocco|
+----------------------------------+-------+--------+--------+
only showing top 20 rows
```

Los 10 hospedajes más baratos

```
+---------------------------------------------------+-----+
|Name                                               |Price|
+---------------------------------------------------+-----+
|Campsite in Rif, Iceland                           |16.0 |
|Campsite in Rif, Iceland                           |16.0 |
|Hotel in Nha Trang, Vietnam                        |17.0 |
|Hut in Guarne, Colombia                            |18.0 |
|Apartment in Bangkok, Thailand                     |19.0 |
|Home in Jumilhac-le-Grand, France                  |20.0 |
|Earthen home in Mueang Chiang Mai District, Thailand|22.0 |
|Earthen home in Mueang Chiang Mai District, Thailand|22.0 |
|Condo in Melaka, Malaysia                          |23.0 |
|Place to stay in Thành phố Hội An, Vietnam         |24.0 |
+---------------------------------------------------+-----+
only showing top 10 rows
```

Filtro de hospedajes por país únicamente en Colombia

```
+-------------------------------+-----+--------+--------+
|Name                           |Price|Category|Country |
+-------------------------------+-----+--------+--------+
|Tiny home in Medellín, Colombia|67.0 |Standard| Colombia|
|Cabin in Medellín, Colombia    |51.0 |Standard| Colombia|
|Place to stay in Retiro, Colombia|107.0|Standard| Colombia|
|Chalet in Cajicá, Colombia     |40.0 |Standard| Colombia|
|Cabin in El Peñol, Colombia    |220.0|Standard| Colombia|
|Cabin in Santa Marta, Colombia |61.0 |Standard| Colombia|
|Hut in Guarne, Colombia        |18.0 |Standard| Colombia|
+-------------------------------+-----+--------+--------+
```

Orden de acuerdo al rating o calificación del hospedaje

```
+-----------------+-----------------+---------+-----+-----------------+-------------+--------+-----+----------+---------+
|             Name|           Detail|     Date|Price|Review and rating|Number of bed|Category|Month|      city|  country|
+-----------------+-----------------+---------+-----+-----------------+-------------+--------+-----+----------+---------+
| Home in Los Angeles|Hollywood Sunset ...|May 13 - 18|272.0|         5.0 (99)|  2 king beds|Standard|  May|      NULL|     NULL|
|Villa in Kathu, T...|Tiny Poolvilla in...|Jul 25 - 30|119.0|         5.0 (98)|   1 king bed|Standard|  Jul|  Thailand| Thailand|
|Villa in Kathu, T...|Tiny Poolvilla in...|Jul 25 - 30|119.0|         5.0 (98)|   1 king bed|Standard|  Jul|  Thailand| Thailand|
|Villa in Morongo ...|Pink Galaxy | Pri...|Jun 7 - 12|283.0|          5.0 (9)|  1 queen bed|Standard|  Jun|California|       US|
|Farm stay in Şark...|Off-Grid vineyard...| Jun 7 - 13|144.0|          5.0 (9)|       3 beds|Standard|  Jun|    Turkey|   Turkey|
|Cottage in Boscas...|A Boscastle barn,...|Jun 11 - 16|141.0|          5.0 (9)|   1 king bed|Standard|  Jun|        UK|       UK|
|Room in Wien, Aus...|Central luxury /H...|Sep 17 - 22| 92.0|          5.0 (9)|       3 beds|Standard|  Sep|   Austria|  Austria|
|Room in Kecamatan...|Odesa Villa1 Ubud...|Jun 11 - 16| 29.0|          5.0 (9)|        1 bed|Standard|  Jun| Indonesia|Indonesia|
|   Home in Ellington|  Crystal Lake House| May 1 - 6|155.0|          5.0 (9)|       4 beds|Standard|  May|      NULL|     NULL|
|  Tiny home in Greer|*BRAND NEW* Tiny ...| May 1 - 6| 90.0|         5.0 (86)|       2 beds|Standard|  May|      NULL|     NULL|
+-----------------+-----------------+---------+-----+-----------------+-------------+--------+-----+----------+---------+
only showing top 10 rows
```

Extracción de la valoración

```
+------------------+----------------+------------+-----+-----------------+-----------+--------+--------+-----------+----------+----------+
|              Name|          Detail|        Date|Price|Review and rating|Number of bed|Category|Month|       city|   country|valoración|
+------------------+----------------+------------+-----+-----------------+-----------+--------+--------+-----------+----------+----------+
| Home in Los Angeles|Hollywood Sunset ...|May 13 - 18|272.0|          5.0 (99)| 2 king beds|Standard|   May|       NULL|      NULL|       5.0|
|Villa in Kathu, T...|Tiny Poolvilla in...|Jul 25 - 30|119.0|          5.0 (98)|  1 king bed|Standard|   Jul|   Thailand|  Thailand|       5.0|
|Villa in Kathu, T...|Tiny Poolvilla in...|Jul 25 - 30|119.0|          5.0 (98)|  1 king bed|Standard|   Jul|   Thailand|  Thailand|       5.0|
|Villa in Morongo ...|Pink Galaxy | Pri...|Jun 7 - 12|283.0|           5.0 (9)| 1 queen bed|Standard|   Jun| California|        US|       5.0|
|Farm stay in Şark...|Off-Grid vineyard...|Jun 7 - 13|144.0|           5.0 (9)|      3 beds|Standard|   Jun|     Turkey|    Turkey|       5.0|
|Cottage in Boscas...|A Boscastle barn,...|Jun 11 - 16|141.0|           5.0 (9)|  1 king bed|Standard|   Jun|         UK|        UK|       5.0|
|Room in Wien, Aus...|Central luxury /H...|Sep 17 - 22| 92.0|           5.0 (9)|      3 beds|Standard|   Sep|    Austria|   Austria|       5.0|
|Room in Kecamatan...|Odesa Villa1 Ubud...|Jun 11 - 16| 29.0|           5.0 (9)|       1 bed|Standard|   Jun|  Indonesia| Indonesia|       5.0|
|   Home in Ellington|  Crystal Lake House| May 1 - 6|155.0|           5.0 (9)|      4 beds|Standard|   May|       NULL|      NULL|       5.0|
|  Tiny home in Greer|*BRAND NEW* Tiny ...| May 1 - 6| 90.0|          5.0 (86)|      2 beds|Standard|   May|       NULL|      NULL|       5.0|
+------------------+----------------+------------+-----+-----------------+-----------+--------+--------+-----------+----------+----------+
only showing top 10 rows
```

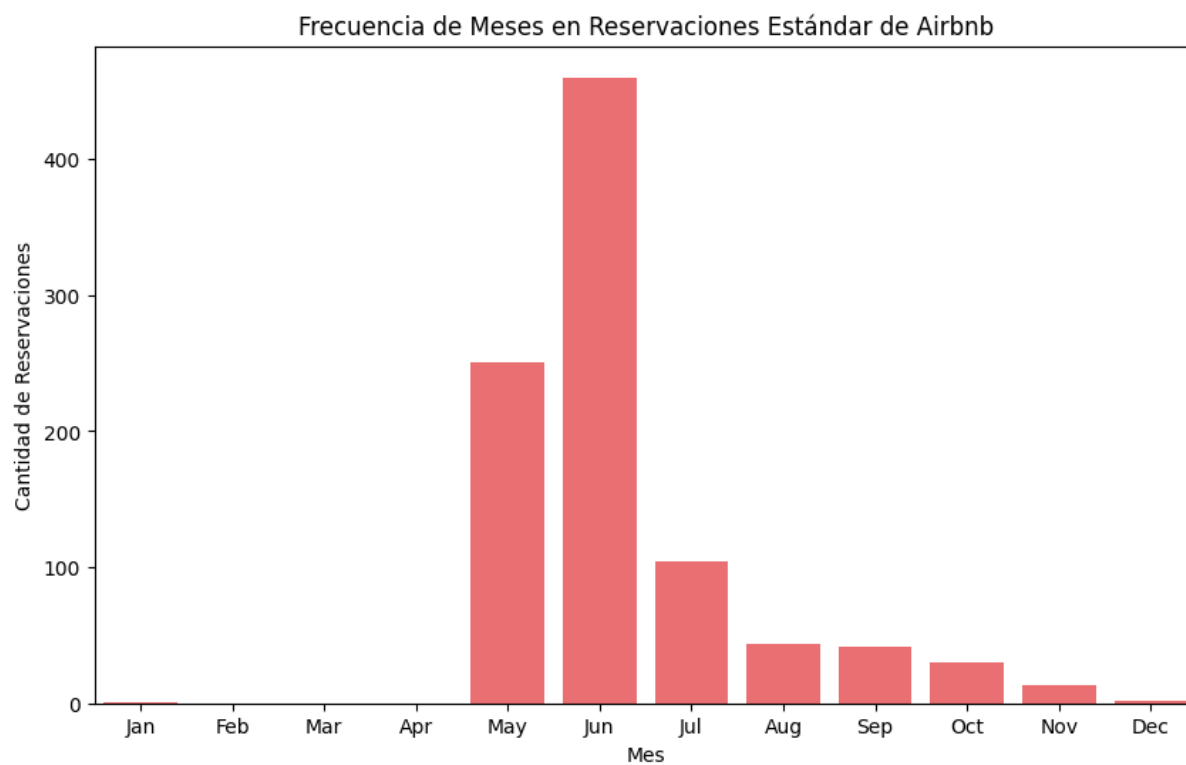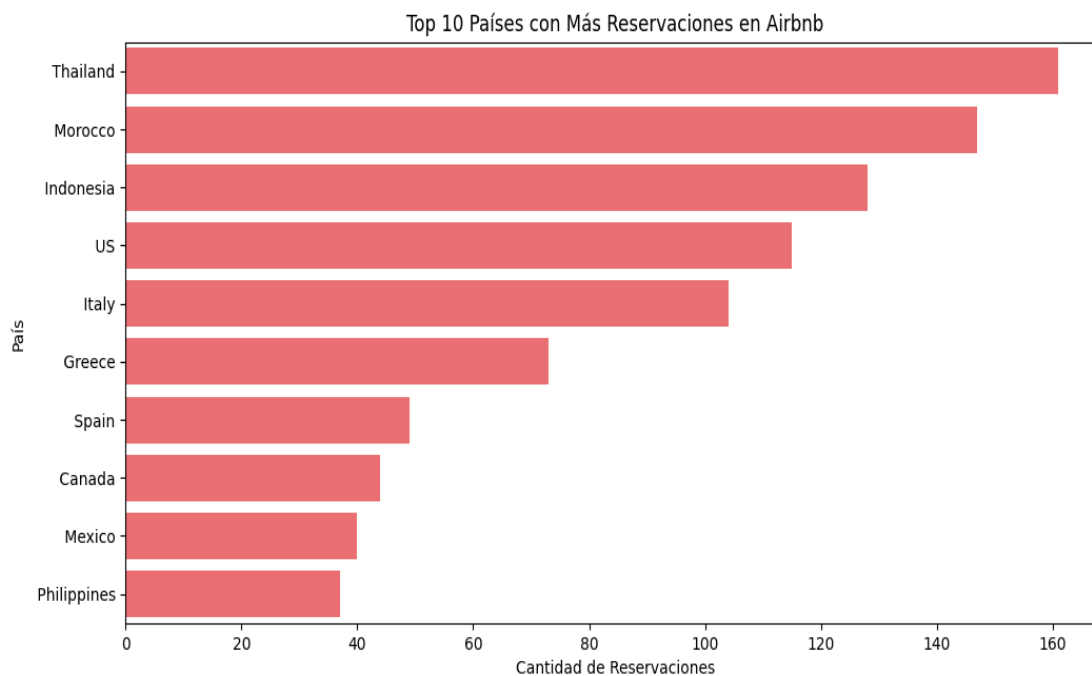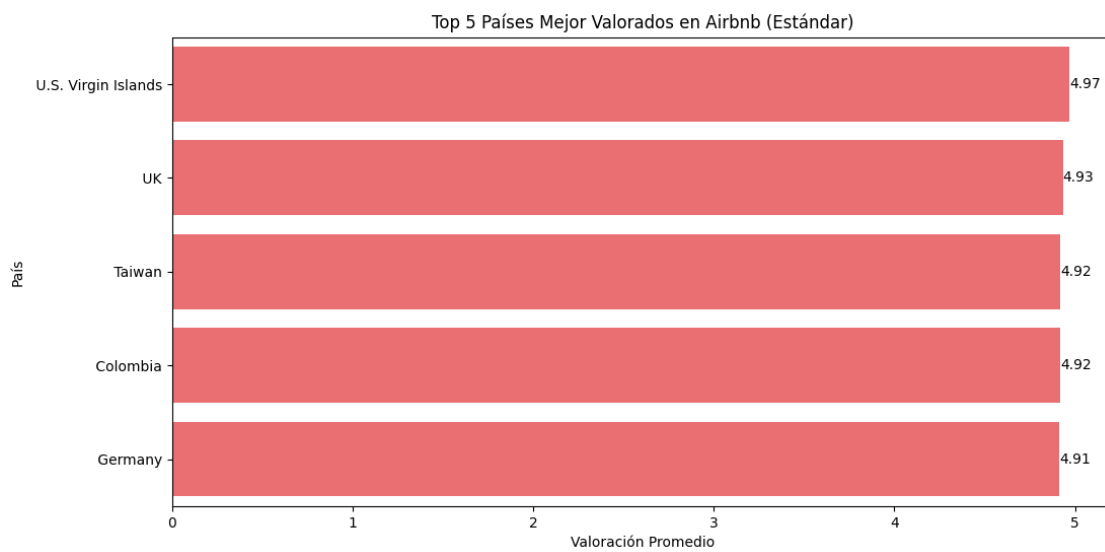Visualización de datos:



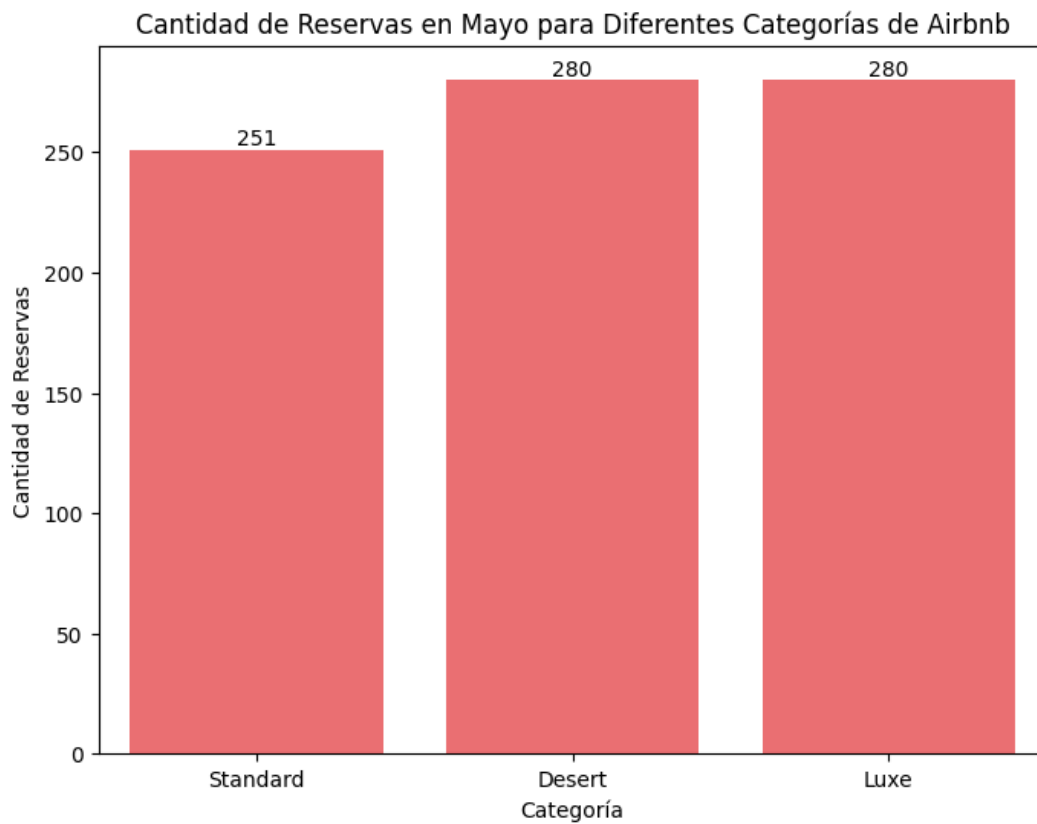Cantidad de reservas por País en Airbnb -Standar



Cantidad de reservas por País en Airbnb - Desierto



Cantidad de reservas por País en Airbnb - Luxe

Top 10 Países con Más Reservaciones en Airbnb



Frecuencia de Meses en Reservaciones Estándar de Airbnb

**Cantidad de Reservas en Mayo para Diferentes Categorías de Airbnb**



**Top 5 Países Mejor Valorados en Airbnb (Estándar)**

Top 10 Lugares Más Caros en Airbnb
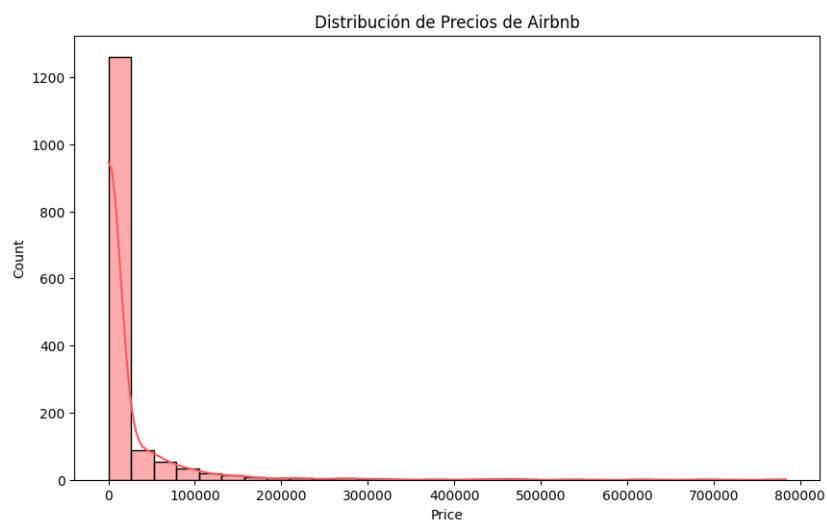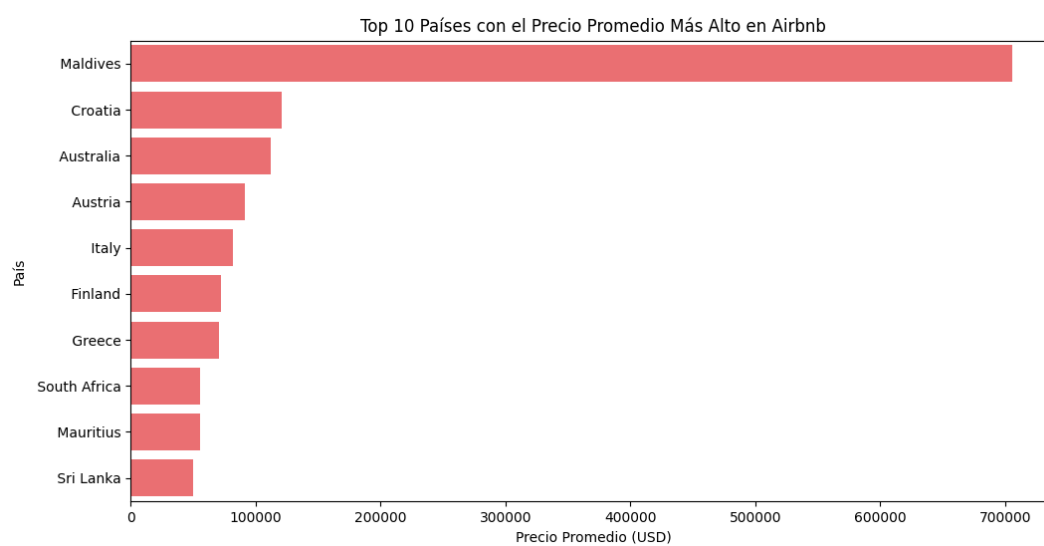


Top 10 Países con el Precio Promedio Más Alto en Airbnb



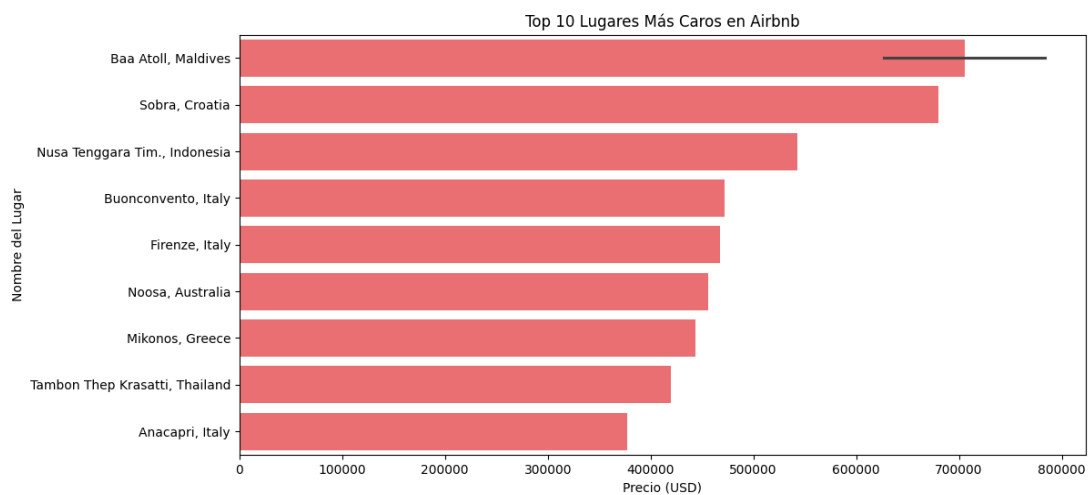Distribución de Precios de Airbnb

## 3. Referencias

Joy Shil. (2023). Airbnb Listing Data 2023 [Data set]. Kaggle. https://doi.org/10.34740/KAGGLE/DSV/5793330