

# Informe y presentación

Teresa Álvarez de Portugal

May 2023

## 1 Introduction

Aquí transmitiré los resultados de los análisis realizados. Voy a explicar paso a paso cada uno de los apartados con las conclusiones correspondientes de las tareas que he realizado. Incluiré principalmente los resultados de las tareas realizadas sobre los datos obtenidos a través de la ejecución del código contenido en el documento técnico.

### 1.1 Ejercicio 1

Esta tarea consiste en conocer con qué tipo de datos contamos. Para ello, he tenido que categorizarlos según su estructura y presentarlos de la siguiente manera:

NOMBRE DEL CAMPO 1	TIPO DE DATO
Activity Period	Integer
Operating Airline 4	Object
Operating Airline IATA Code	Object
Published Airline	Object
Published Airline IATA Code	Object
GEO Summary	Object
GEO Region	Object
Activity Type Code	Object
Price Category Code	Object
Terminal	Object
Boarding Area	Object
Passenger Count	Integer
Adjusted Activity Type Code	Object
Adjusted Passenger count	Integer
Year	Integer
Month	Object

### 1.2 Ejercicio 2

Este ejercicio consiste en desarrollar una serie de puntos que incluyen preguntas que son las que contestaré aquí.

-¿Cuántas compañías diferentes aparecen en el fichero?

Aparecen un total de 77 compañías. Para poder contestar a esta pregunta combiné dos columnas que contenían los nombres de las compañías y saqué los valores únicos.

-¿Cuántos pasajeros tiene de media los vuelos de cada compañía?

Como en el caso anterior, había dos columnas que contenían número de pasajeros, pero en este caso las comparé y elegí una de ellas. Utilicé la lista en la que había guardado las compañías

```
array([[ 'ATA Airlines', 'ATA Airlines', 8744.6363636364],
 [ 'Aer Lingus', 'Aer Lingus', 4407.183673469388],
 [ 'Aeromexico', 'Aeromexico', 5463.622222222222],
 [ 'Air Berlin', 'Air Berlin', 2320.75],
 [ 'Air Canada ', 'Air Canada ', 18251.560109289618],
 [ 'Air Canada Jazz', 'Air Canada ', 294.2142857142857],
 [ 'Air China', 'Air China', 6618.335907335907],
 [ 'Air France', 'Air France', 11589.097510390905],
 [ 'Air India limited', 'Air India limited', 2834.5],
 [ 'Air New Zealand', 'Air New Zealand', 7452.339768339768],
 [ 'Airtran Airways', 'Airtran Airways', 10569.238938053897],
 [ 'Alaska Airlines', 'Alaska Airlines', 17251.637816245006],
 [ 'All Nippon Airways', 'All Nippon Airways', 6385.523255813953],
 [ 'Allegiant Air', 'Allegiant Air', 1516.612],
 [ 'American Airlines', 'American Airlines', 127164.38970588235],
 [ 'American Eagle Airlines', 'American Airlines',
 4006.5283018867926],
 [ 'Aeriflight', 'Aeriflight', 5.0],
 [ 'Asiana Airlines', 'Asiana Airlines', 5902.961240310077],
 [ 'Atlantic Southeast Airlines', 'Delta Air Lines',
 2176.909090909091],
 [ 'Atlas Air, Inc', 'Atlas Air, Inc', 34.0],
 [ 'BelAir Airlines', 'BelAir Airlines', 415.36363636364],
 [ 'Boeing Company', 'Boeing Company', 18.0],
 [ 'British Airways', 'British Airways', 17625.124031007752],
 ...
 [ 'Virgin Atlantic', 'Virgin Atlantic', 9847.10465116279],
 [ 'WestJet Airlines', 'WestJet Airlines', 5338.155339805825],
 [ 'World Airways', 'World Airways', 261.6666666666667],
 [ 'Xl Airways France', 'Xl Airways France', 2223.1612803225805],
 [ 'Xtra Airways', 'Xtra Airways', 73.0]], dtype=object)
```

anteriormente, filtré el dataframe para las aerolíneas presentes en esa lista y calculé el número medio de pasajeros por compañía. El resultado es el siguiente:

Como se aprecia son varios datos debido a que son 77 compañías, por lo que en el notebook se puede comprobar la respuesta mejor.

-Eliminar los registros duplicados por el campo GEO Region, manteniendo únicamente aquel con mayor número de pasajeros

Para realizar esta tarea agrupé los datos de la columna pedida. Después calculé el valor máximo de las columnas de pasajeros para cada grupo de GEO Region. Y reestablecí el índice porque después me pide que lo guarde en un dataset nuevo. Los resultados se pueden observar bien en la siguiente imagen.

	GEO Region	Passenger Count	Adjusted Passenger Count
0	US	659837	659837
1	Canada	39798	39798
2	Asia	86398	86398
3	Europe	48136	48136
4	Australia / Oceania	12973	12973
5	Mexico	29206	29206
6	Central America	8970	8970
7	Middle East	14769	14769
8	South America	3685	3685

### 1.3 Ejercicio 3

En este ejercicio se nos pide realizar un análisis descriptivo de los datos usando Dask y calculando la media y la desviación estándar de cada elemento del conjunto de datos. Para esto primero observé que solo había cuatro columnas con valores numéricos, por lo que entendí que tenía que codificar el dataframe para poder trabajar comodamente. Elegí entre la técnica de codificación de Label Encoding y One-Hot Encoding, porque ambas son usadas comunmente en el procesamiento de datos antes de realizar el análisis. Para este caso que tenía que calcular tanto media como desviación estándar usé Label Encoding, el cual asigna valores numéricos únicos a cada categoría, lo que haría mas fácil el cálculo.

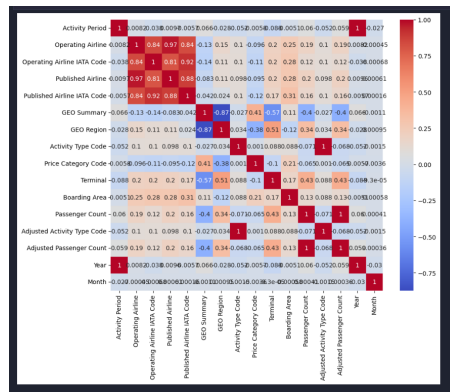
-Mis conclusiones fueron las siguientes:

A parte de calcular la media y la desviación estandar, he calculado la moda porque creo que es mas relevante en el caso de analizar algunas columnas.

La media representa el valor promedio de los datos de una variable. En este caso sería efectiva sobre las columnas como la de 'Passenger count', 'ADjusted passenger count', 'Price code category', 'GEO summary'. En las demas son datos en los que necesariamente los valores no tienen por que tener un orden, con lo que sería mas efectivo usar la moda para ver que valor es el que mas se repite.

En cuanto a la desviación, es una medida de dispersión que indica cuánto se alejan los datos de la media. Los que tienen desviación estandar alta quiere decir que tienen los valores más dispersos y hay mayor variabilidad. Donde es baja significa que los valores están más cercanos a la media. Para interpretar todo esto también hay que tener en cuenta la cantidad de observaciones. En las columnas donde hay pasajeros por ejemplo, hay una desviación muy alta porque los valores van desde 2 hasta un número muy alto. En cambio en columnas como el precio o la actividad hay una desviación baja porque hay pocos valores por lo que estarán más cerca de la media.

Después nos pide también realizar un análisis de la correlación cuyo resultado debía ser una matriz de correlación de datos que represente de qué manera están relacionadas las diferentes variables. Como ya tenía el dataframe codificado continúe programando la matriz de correlación que me daba como resultado lo siguiente:



-De donde saqué las siguientes conclusiones:

Las celdas en las que hay correlación alta, en este caso de 0,8 para arriba, son las que tienen una correlación positiva muy fuerte. En este caso se puede observar que las 4 columnas que hablan sobre las aerolíneas y sus códigos IATA están altamente correlacionadas ya que hablan de los mismos datos. Lo mismo ocurre con las columnas 'Activity Type Code' con 'Adjusted Activity Type Code' y 'Passenger Count' con 'Adjusted Passenger count'. De la misma manera ocurre con la correlación entre 'Activity Period' y 'Year' ya que la primera de las mencionadas contiene tanto el año como el mes pero unido y por esa razón no reconoce la correlación con la columna del mes, aunque se puede apreciar cuando se compara.

Donde hay correlación negativa significa que las variables tienden a moverse en direcciones opuestas. Entre las columnas GEO Summary y Geo region es donde está más cerca de 1 por lo que son las más correlacionadas negativamente. La siguiente con estas características serían GEO Summary con Terminal pero no es tan fuerte.

En el resto se puede observar correlación nula o cerca de ella lo que significa que no hay una relación lineal clara entre las variables. Esto no significa que no haya una relación entre las variables, ya que podría haber relaciones no lineales.

Por último, este ejercicio nos pedía seleccionar uno de los algoritmos vistos durante el curso y aplicarlo para después explicar los resultados obtenidos.

Aunque no he podido finalizar la aplicación del algoritmo, mi decisión fue que la mejor opción era usar Apache Spark como framework de procesamiento paralelo y distribuido. En mi opinión es una buena opción gracias a la capacidad que tiene para procesar grandes volúmenes de datos de manera escalable. Esto significa que puede servir bien con el dataset de grandes dimensiones con el que estábamos trabajando. También lo seleccioné porque proporciona una amplia gama de bibliotecas útiles para aplicar algoritmos y análisis específicos al dataset y porque se integra con python a través de la biblioteca PySpark.