# GROUP- 7

**Members:**
**Ritesh Sengar**
**Shravan Honade**
**Nikhil Patil**
**Dhananjay Ghate**

# EDA

## Avg. Distance and Avg. Airtime by Carriers
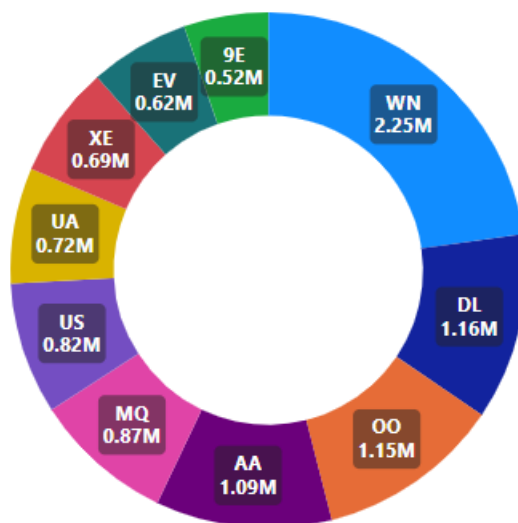
● Average of DISTANCE ● Average of AIR_TIME



Here are some additional observations from the graph:

- Low-cost carriers (LCCs) such as Southwest Airlines (WN) and Spirit Airlines (NK) tend to have shorter average distances and shorter average airtimes than traditional carriers such as United Airlines (UA) and Delta Air Lines (DL). This is because LCCs typically operate shorter flights.

- Regional airlines such as Envoy Air (EV) and Piedmont Airlines (PI) tend to have shorter average distances and shorter average airtimes than major airlines. This is because regional airlines typically operate shorter flights to smaller airports.
- Charter airlines such as NetJets (XE) and Flexjet (HA) tend to have longer average distances and longer average airtimes than commercial airlines. This is because charter airlines typically operate flights to and from smaller airports and private airstrips.

## Top 10 Carriers by Amount of flights

OP_CARRIER  ● WN  ● DL  ● OO  ● AA  ● MQ  ● US  ● UA  ● XE  ● EV  ● 9E



The pie chart shows the top 10 carriers by number of flights in 2017. The carriers are listed in order of decreasing number of flights, with American Airlines (AA) at the top and EV9E at the bottom.

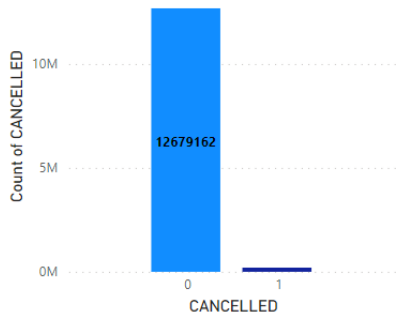The top 5 carriers by number of flights are:

1. American Airlines (AA)
2. Delta Air Lines (DL)
3. Ryanair Group (RYR)
4. United Airlines (UA)
5. Southwest Airlines (WN)

These carriers account for over 60% of all flights in 2017.

## Total Flights cancelled from Top 10 originating Places:

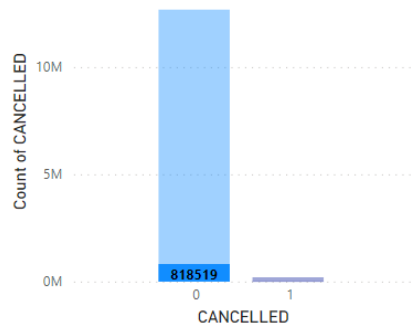| CANCELLATION_CODE | Count of CANCELLATION_CODE ▼ |
|---|---|
| | 12679162 |
| B | 93158 |
| A | 77724 |
| C | 29352 |
| D | 59 |
| **Total** | **12879455** |

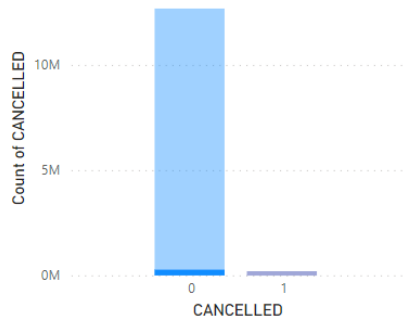CANCELLED ● 0 ● 1



Top 10 places from which flights originated

## Maximum number of Flights cancelled from Originating Location:

| CANCELLATION_CODE | Count of CANCELLATION_CODE ▾ |
|---|---|
|  | 818519 |
| A | 5737 |
| B | 5625 |
| C | 2231 |
| D | 3 |
| **Total** | **832115** |

CANCELLED ● 0 ● 1



Top 10 places from which flights originated

## Minimum number of flights cancelled from Originating Location:

| CANCELLATION_CODE | Count of CANCELLATION_CODE ▾ |
|---|---|
|  | 272251 |
| A | 1709 |
| B | 1264 |
| C | 1264 |
| D | 2 |
| **Total** | **276490** |

CANCELLED ● 0 ● 1



Top 10 places from which flights originated

# PySpark Application deployment on Kubernetes:

**Artifact Registry**

| | |
|---|---|
| ☰ | Repositories |
| ⚙ | Settings |

← Images for k8-...    🗑 DELETE    ⋮    C

📁 us-east1-docker.pkg.dev  ›  📁 dataproc-pyspark-404719  ›  ■ k8-pyspark-docker

## Repository Details

| | |
|---|---|
| Format | Docker |
| Type | Standard |

⌄ SHOW MORE

≡ Filter    Enter property name or value                                    ❓

| ☐ | Name ↑ | Created | Updated |
|---|---|---|---|
| ☐ | 🐳 myk8spark | Nov 24, 2023 | 13 days ago |
| ☐ | 🐳 myk8spark/myk8spark | 13 days ago | 13 days ago |
| ☐ | 🐳 pysparkml | 1 hour ago | Just now |

```yaml
 1   apiVersion.    .../v1beta2"
 2   kind: SparkApplication
 3   metadata:
 4     name: pyspark-ml
 5     namespace: default
 6     labels:
 7       app: pyspark-ml
 8   spec:
 9     type: Python
10     pythonVersion: "3"
11     mode: cluster
12     image: "us-east1-docker.pkg.dev/dataproc-pyspark-404719/k8-pyspark-docker/pysparkml"
13     # https://skaffold.dev/docs/environment/local-cluster/
14     #  Skaffold's direct loading of images into a local cluster does mean that resources
15     # imagePullPolicy: Always may fail as the images are not be pushed to the remote reg
16     # On Docker for Desktop, don't specify imagePullPolicy
17     imagePullPolicy: Always
18     mainApplicationFile: local:///opt/spark/work-dir/FlightCancellationPrediction.py
19     sparkConf:
20       "spark.ui.port": "4040"
21     sparkVersion: "3.2"
22     restartPolicy:
23       type: Never
24     driver:
25       coreLimit: "1"
26       coreRequest: "1m"
27       memory: "512m"
28       labels:
29         version: "3.2"
30       serviceAccount: spark
31     executor:
32       coreLimit: "1"
33       coreRequest: "1m"
34       instances: 1
35       memory: "512m"
36       labels:
```