

DEVELOPING A FLIGHT DELAY MODEL USING MACHINE LEARNING

TEAM ID: PNT2022TMID27775

Analyze The Data

- How the information is stored in a DataFrame or Python object affects what we can do with it and the outputs of calculations as well. There are two main types of data : numeric and text data types.
- Numeric data types include integers and floats.
- Text data type is known as Strings in Python, or Objects in Pandas. Strings can contain numbers and / or characters.
- For example, a string might be a word, a sentence, or several sentences.
- Will see how our dataset is, by using info() method.

```
dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11231 entries, 0 to 11230
Data columns (total 26 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   YEAR                                  11231 non-null  int64
1   QUARTER                              11231 non-null  int64
2   MONTH                                11231 non-null  int64
3   DAY_OF_MONTH                         11231 non-null  int64
4   DAY_OF_WEEK                          11231 non-null  int64
5   UNIQUE_CARRIER                     11231 non-null  object
6   TAIL_NUM                             11231 non-null  object
7   FL_NUM                               11231 non-null  int64
8   ORIGIN_AIRPORT_ID                   11231 non-null  int64
9   ORIGIN                               11231 non-null  object
10  DEST_AIRPORT_ID                     11231 non-null  int64
11  DEST                                 11231 non-null  object
12  CRS_DEP_TIME                        11231 non-null  int64
13  DEP_TIME                            11124 non-null  float64
14  DEP_DELAY                           11124 non-null  float64
15  DEP_DEL15                           11124 non-null  float64
16  CRS_ARR_TIME                        11231 non-null  int64
17  ARR_TIME                            11116 non-null  float64
18  ARR_DELAY                           11043 non-null  float64
19  ARR_DEL15                           11043 non-null  float64
20  CANCELLED                           11231 non-null  float64
21  DIVERTED                            11231 non-null  float64
22  CRS_ELAPSED_TIME                    11231 non-null  float64
23  ACTUAL_ELAPSED_TIME                 11043 non-null  float64
24  DISTANCE                            11231 non-null  float64
25  Unnamed: 25                         0 non-null     float64
dtypes: float64(12), int64(10), object(4)
memory usage: 2.2+ MB
```

- As you can see in our dataset both numerical and categorical data are present, but it is not necessary that all the continuous data which we are seeing has to be continuous in nature. There may be a case that some categorical data is in the form of numbers but when we perform info() operation we will get numerical output. So, we need to take care of those types of data also.

```
dataset.describe()
```

	YEAR	QUARTER	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	FL_NUM	ORIGIN_AIRPORT_ID	DEST_AIRPORT_ID	CRS_DEP_TIME	DEP_
count	11231.0	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11124.00
mean	2016.0	2.544475	6.628973	15.790758	3.960199	1334.325617	12334.516695	12302.274508	1320.798326	1327.11
std	0.0	1.090701	3.354678	8.782056	1.995257	811.875227	1595.026510	1601.988550	490.737845	500.31
min	2016.0	1.000000	1.000000	1.000000	1.000000	7.000000	10397.000000	10397.000000	10.000000	1.00
25%	2016.0	2.000000	4.000000	8.000000	2.000000	624.000000	10397.000000	10397.000000	905.000000	905.00
50%	2016.0	3.000000	7.000000	16.000000	4.000000	1267.000000	12478.000000	12478.000000	1320.000000	1324.00
75%	2016.0	3.000000	9.000000	23.000000	6.000000	2032.000000	13487.000000	13487.000000	1735.000000	1739.00
max	2016.0	4.000000	12.000000	31.000000	7.000000	2853.000000	14747.000000	14747.000000	2359.000000	2400.00

8 rows × 22 columns

- Describe() functions are used to compute values like count, mean, standard deviation give a summary type of data