

# Welcome to Safran Lab 1

Every day, more than 80,000 commercial flights take place around the world, operated by hundreds of airlines. For all aircraft take-off weight exceeding 27 tons, a regulatory constraint requires companies to systematically record and analyse all flight data, for the purpose of improving the safety of flights. Flight Data Monitoring strives to detect and prioritize deviations from standards set by the aircraft manufacturers, the authorities of civil aviation in the country, or even companies themselves. Such deviations, called events, are used to populate a database that enables companies to identify and monitor the risks inherent to these operations.

This notebook is designed to let you manipulate real aeronautical data, provided by the Safran Group. It is divided in two parts: the first part deals with the processing of raw data, you will be asked to visualize the data, understand what variables require processing and perform the processing for some of these variables. The second part deals with actual data analysis, and covers some interesting problems. We hope to give you some insights of the data scientist job and give you interesting and challenging questions.

## Part 1: Data processing

### Loading raw data

#### Context

You will be provided with 780 flight records. Each is a full record of a flight starting at the beginning of the taxi out phase and terminating at the end of the taxi in phase. The sample rate is 1 Hz. Please be aware that due to side effects the very beginning of the record may be faulty. This is something to keep in mind when we will analyse the data.

Each flight data is a collection of time series resumed in a dataframe, the columns variables are described in the schema below:

name	description	unit
TIME	elapsed seconds	second
LATP_1	Latitude	degree °
LONP_1	Longitude	degree °
RALT1	Radio Altitude, sensor 1	feet
RALT2	Radio Altitude, sensor 2	feet
RALT3	Radio Altitude, sensor 3	feet
ALT_STD	Relative Altitude	feet
HEAD	head	degree °
PITCH	pitch	degree °
ROLL	roll	degree °

IAS	Indicated Air Speed	m/s
N11	speed N1 of the first engine	%
N21	speed N2 of the first engine	%
N12	speed N1 of the second engine	%
N22	speed N2 of the second engine	%
AIR_GROUND	1: ground, 0: air	boolean

**Note :** TIME represents the elapsed seconds from today midnight. You are not provided with an absolute time variable that would tell you the date and hour of the flights.

**Acquire expertise about aviation data**

You will need some expertise about the signification of the variables. Latitude and longitude are quite straightforward. Head, Pitch and Roll are standards orientation angles, check this [image](#) to be sure. RALT\* are coming from three different radio altimeters, they measure the same thing but have a lot of missing values and are valid only under a threshold altitude (around 5000 feet). Alt\_std is the altitude measured from the pressure (it basically comes from a barometer), it is way less accurate that a radio altimeter but provides values for all altitudes. N1\* and N2\* are the rotational speeds of the engine sections expressed as a percentage of a nominal value. Some good links to check out to go deeper:

- [about phases of flight](#)
- [pitch-roll-head](#)
- [about N\\*\\* variables I](#)
- [about N\\*\\* variables II](#)
- [how altimeters work](#)
- [about runway naming](#)

```
In [161]: # Set up

BASE_DIR = "/mnt/safran/TP1/data/"

from os import listdir
from os.path import isfile, join

import glob

import matplotlib as mpl
mpl.rcParams["axes.grid"] = True
import matplotlib.pyplot as plt
%matplotlib inline

import numpy as np
import pandas as pd
pd.options.display.max_columns = 50

from datetime import datetime

from haversine import haversine

def load_data_from_directory(DATA_PATH, num_flights):
    files_list = glob.glob(join(DATA_PATH, "*pk1"))
    print("There are %d files in total" % len(files_list))
```

```
files_list = files_list[:num_flights]
print("We process %d files" % num_flights)
dfs = []
p = 0
for idx, f in enumerate(files_list):
    if idx % int(len(files_list)/10) == 0:
        print(str(p*10) + "%: [" + "#" * p + " " * (10-p) + "]", end="\r")
        p += 1
    dfs.append(pd.read_pickle(f))
print(str(p*10) + "%: [" + "#" * p + " " * (10-p) + "]", end="\r")

return dfs
```

Execute the cell below to load the data for part 1

```
In [162]: num_flights = 780
flights = load_data_from_directory(BASE_DIR + "part1/flights", num_fligh
ts)
for f in flights:
    l = len(f)
    new_idx = pd.date_range(start=pd.Timestamp("now").date(), periods=l,
freq="S")
    f.set_index(new_idx, inplace=True)
```

There are 780 files in total  
We process 780 files  
100%: [#####]

The data is loaded with pandas. Please take a look at the [pandas cheat sheet](#) if you have any doubt. You are provided with 780 dataframes, each of them represents the records of the variables defined above during a whole flight.

flights is a list where each item is a dataframe storing the data of one flight. There is no particular ordering in this list. All the flights depart from the same airport and arrive at the same airport. These airports are hidden to you and you will soon understand how.

For example flights[0] is a dataframe, representing one flight.

```
In [3]: # Give an alias to flights[0] for convenience
f = flights[0]
flights[0].head()
```

Out[3]:

	TIME	LATP_1	LONP_1	HEAD	PITCH	ROLL	IAS	RALT1	RALT2	RALT3	ALT
2017-05-31 00:00:00	0	11.7731	33.3118	350.5	-0.9	-1.0	45.0	NaN	NaN	-4.0	-83.0
2017-05-31 00:00:01	1	11.7731	33.3118	350.5	-0.9	-1.0	45.0	-4.0	NaN	NaN	-83.0
2017-05-31 00:00:02	2	50.7330	21.2602	350.5	-0.9	-1.0	45.0	NaN	NaN	-4.0	-83.0
2017-05-											

31 00:00:03	3	48.0318	12.5277	350.1	-0.9	-1.0	45.0	NaN	-5.0	NaN	-84.0
2017-05-31 00:00:04	4	48.0318	12.5277	350.1	-0.9	-1.0	45.0	NaN	NaN	-4.0	-83.0

You can select a column by indexing by its name.

```
In [4]: f["PITCH"].describe()
```

```
Out[4]: count      7704.000000
mean         2.005101
std          2.378875
min         -2.800000
25%          0.500000
50%          1.800000
75%          2.800000
max          16.000000
Name: PITCH, dtype: float64
```

Use `iloc[]` to select by line number, either the whole dataframe to obtain all the variables of a dataframe...

```
In [5]: f.iloc[50:60]
```

```
Out[5]:
```

	TIME	LATP_1	LONP_1	HEAD	PITCH	ROLL	IAS	RALT1	RALT2	RALT3	ALT
2017-05-31 00:00:50	50	48.0314	12.5278	352.9	-0.9	-1.0	45.0	NaN	NaN	-4.0	-82.0
2017-05-31 00:00:51	51	48.0314	12.5278	352.9	-0.9	-1.0	45.0	NaN	-5.0	NaN	-81.0
2017-05-31 00:00:52	52	48.0314	12.5278	352.9	-0.9	-1.0	45.0	NaN	NaN	-4.0	-81.0
2017-05-31 00:00:53	53	48.0314	12.5278	352.9	-0.9	-1.0	45.0	-4.0	NaN	NaN	-81.0
2017-05-31 00:00:54	54	48.0314	12.5278	352.9	-0.9	-1.0	45.0	NaN	NaN	-4.0	-81.0
2017-05-31 00:00:55	55	48.0314	12.5278	352.9	-0.9	-1.0	45.0	NaN	-5.0	NaN	-81.0
2017-05-31 00:00:56	56	48.0314	12.5278	352.9	-0.9	-1.0	45.0	NaN	NaN	-4.0	-81.0

<b>2017-05-31 00:00:57</b>	57	48.0316	12.5278	352.9	-0.9	-1.0	45.0	-4.0	NaN	NaN	-81.0
<b>2017-05-31 00:00:58</b>	58	48.0316	12.5278	352.9	-0.9	-1.0	45.0	NaN	NaN	-4.0	-81.0
<b>2017-05-31 00:00:59</b>	59	48.0316	12.5278	352.9	-0.9	-1.0	45.0	NaN	-5.0	NaN	-82.0

...or an individual series.

```
In [6]: f["PITCH"].iloc[50:60]
```

```
Out[6]: 2017-05-31 00:00:50    -0.9
2017-05-31 00:00:51    -0.9
2017-05-31 00:00:52    -0.9
2017-05-31 00:00:53    -0.9
2017-05-31 00:00:54    -0.9
2017-05-31 00:00:55    -0.9
2017-05-31 00:00:56    -0.9
2017-05-31 00:00:57    -0.9
2017-05-31 00:00:58    -0.9
2017-05-31 00:00:59    -0.9
Freq: S, Name: PITCH, dtype: float64
```

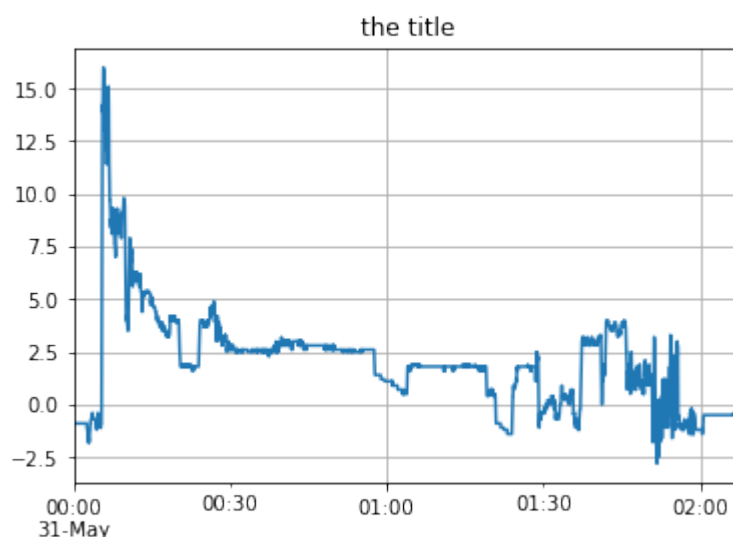
Finally let's work out an example of visualization of a column.

```
In [7]: # Create a figure and one subplot
fig, ax = plt.subplots()

# Give an alias to flights[0] for convenience
f = flights[0]

# Select PITCH column of f and plot the line on ax
f.PITCH.plot(title="the title", ax=ax)
```

```
Out[7]: <matplotlib.axes._subplots.AxesSubplot at 0x7fbd1a9f8710>
```



# Visualization

To perform monitoring of flights, it is necessary to clean up the data. To start, it is important to visualize the data that is available, in order to understand better their properties and the problems associated with them (noise, statistical characteristics, features and other values).

For the following questions do not hesitate to resort to the documentation of pandas for plotting capabilities (for a [dataframe](#) or for a [series](#))

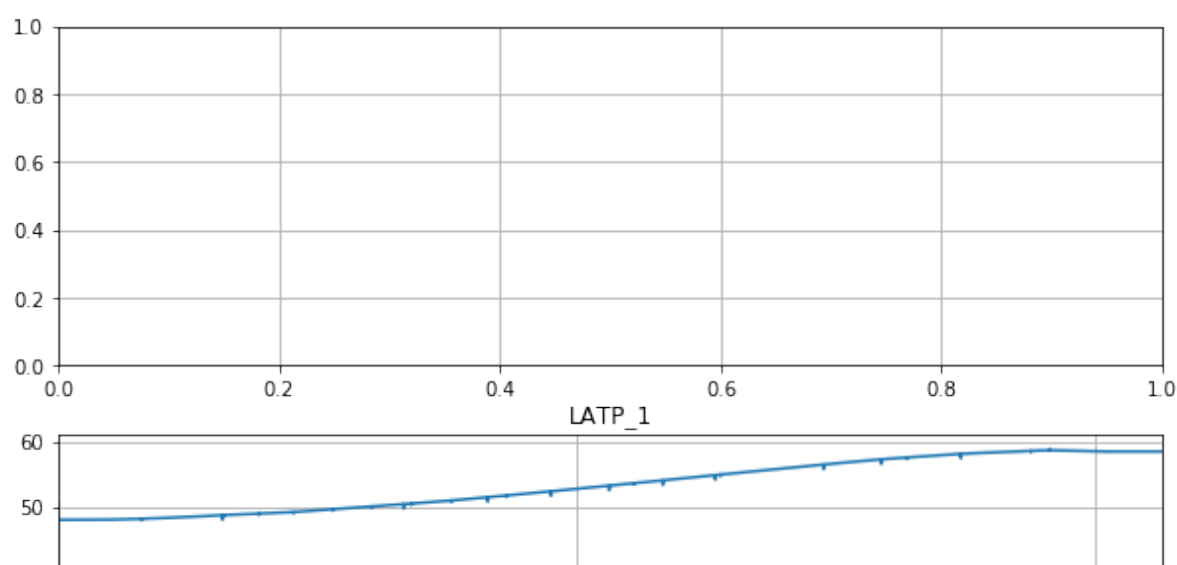
## Question 1 Visualize all the variables

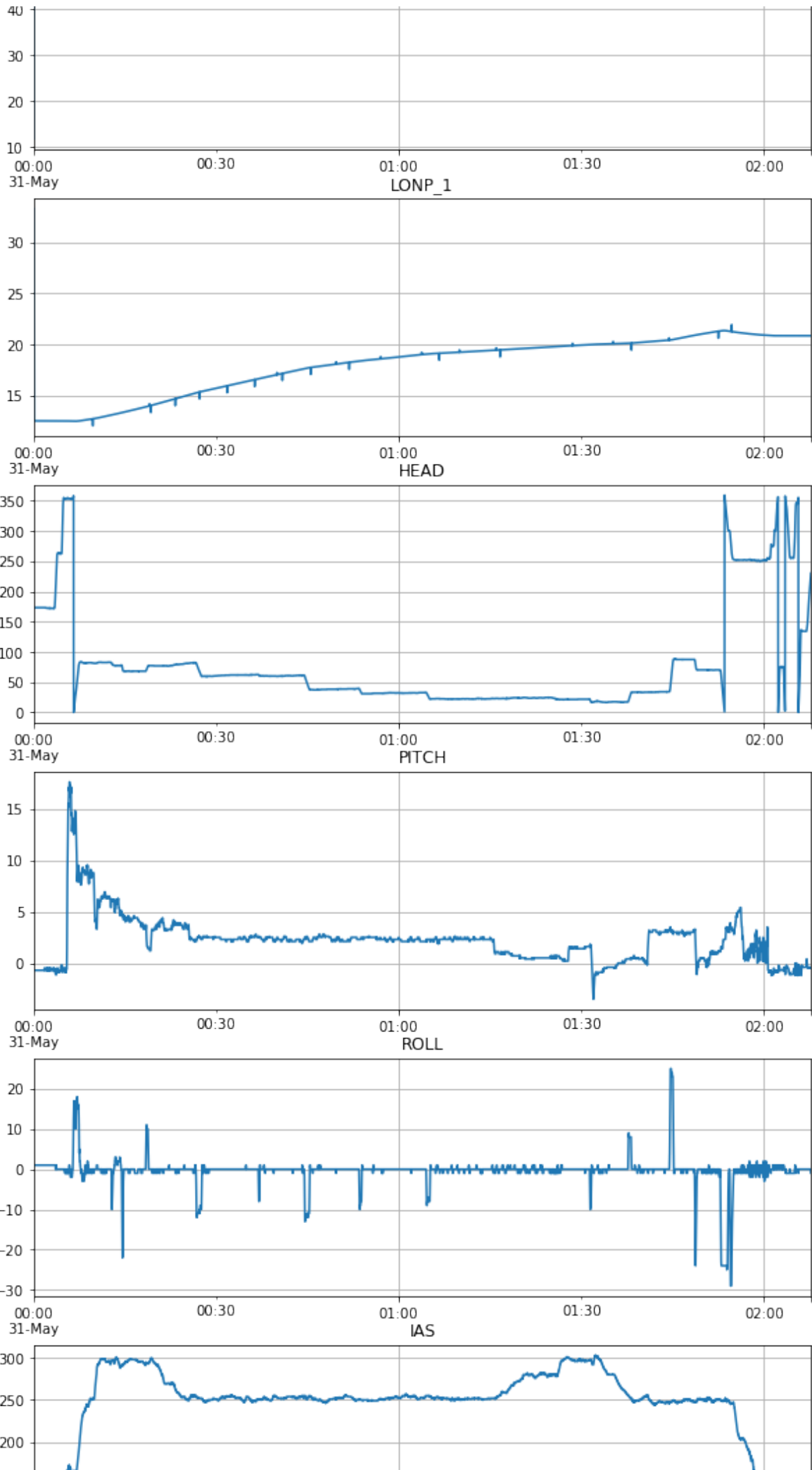
For an arbitrary flight, for example `flights[0]`, visualize all the variables. Would you rather use `plot` or `scatter`? Interpolate the data or not interpolate? Think about NaN values and how they are treated when we plot a series. Comment.

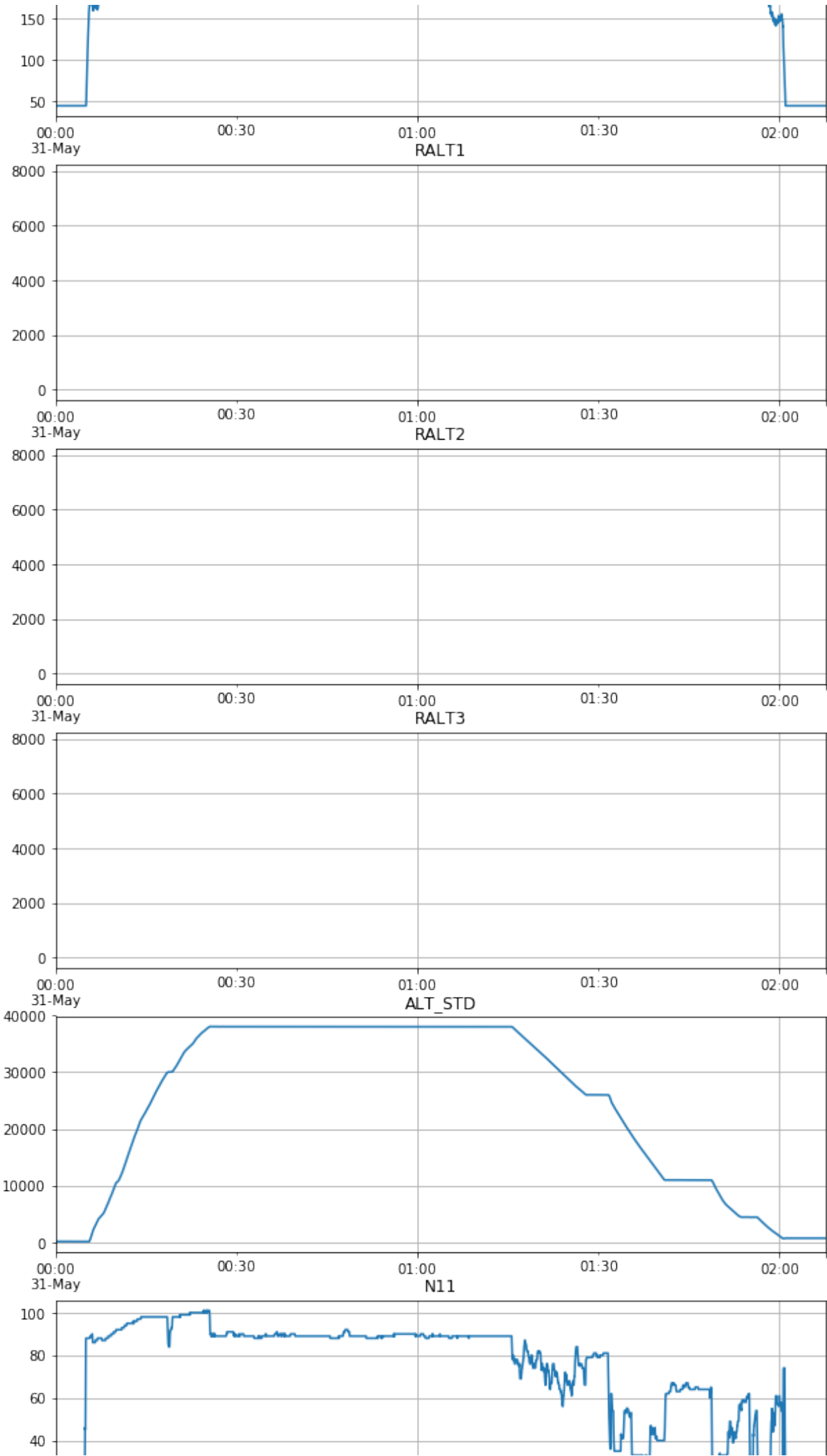
```
In [56]: # LATP_1      LONP_1  HEAD    PITCH  ROLL    IAS      RALT1  RALT2
RALT3  ALT_STD N11      N21      N22      N12      AIR_GROUND
fig, axarr = plt.subplots(16,1, figsize=[10, 60])

f.LATP_1.plot(title='LATP_1', ax=axarr[1])
f.LONP_1.plot(title='LONP_1',ax=axarr[2])
f.HEAD.plot(title='HEAD',ax=axarr[3])
f.PITCH.plot(title='PITCH',ax=axarr[4])
f.ROLL.plot(title='ROLL',ax=axarr[5])
f.IAS.plot(title='IAS',ax=axarr[6])
f.RALT1.plot(title='RALT1',ax=axarr[7])
f.RALT2.plot(title='RALT2',ax=axarr[8])
f.RALT3.plot(title='RALT3',ax=axarr[9])
f.ALT_STD.plot(title='ALT_STD',ax=axarr[10])
f.N11.plot(title='N11',ax=axarr[11])
f.N12.plot(title='N12',ax=axarr[12])
f.N21.plot(title='N21',ax=axarr[13])
f.N22.plot(title='N22',ax=axarr[14])
f.AIR_GROUND.plot(title='AIR_GROUND',ax=axarr[15])
```

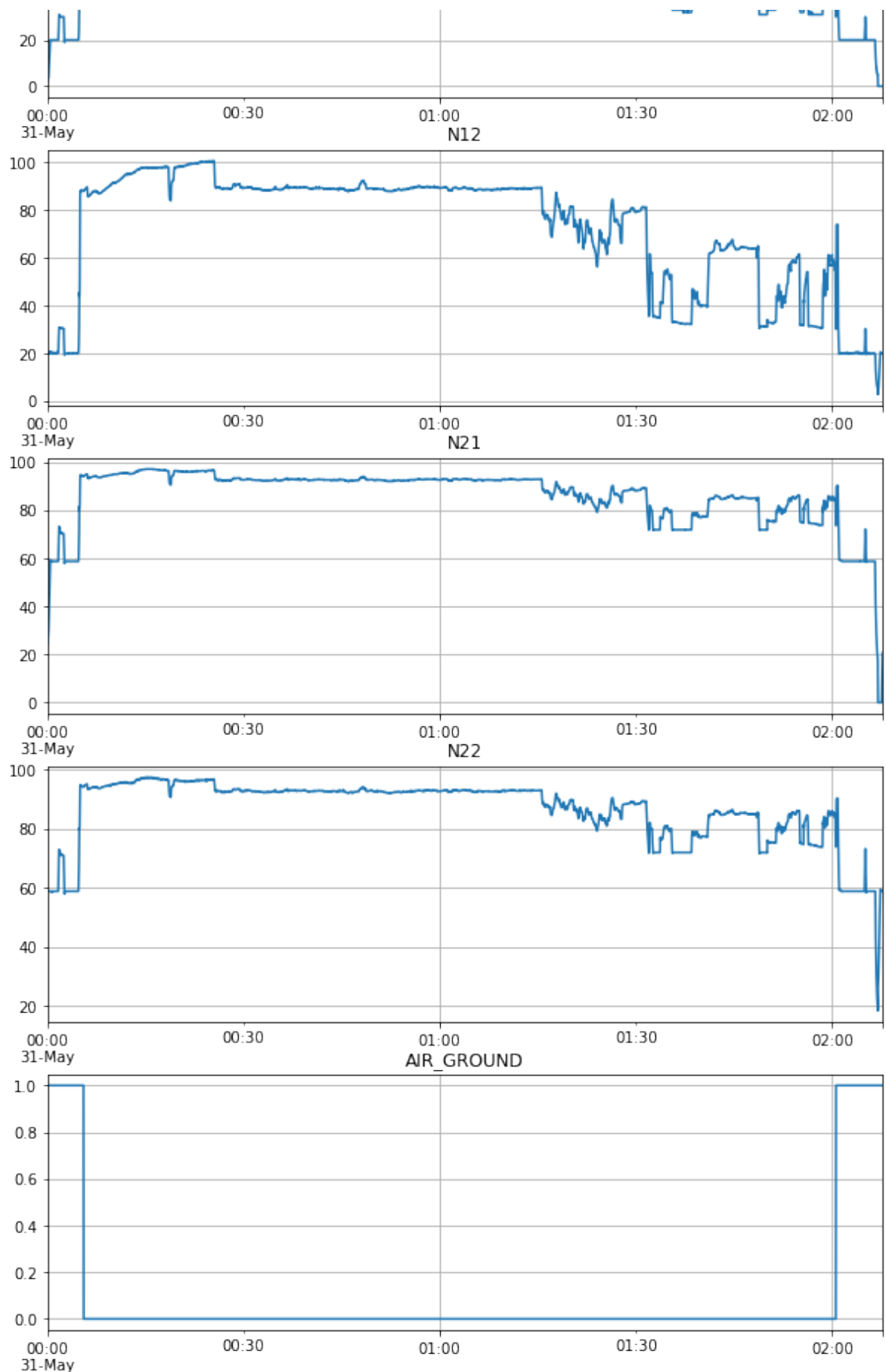
```
Out[56]: <matplotlib.axes._subplots.AxesSubplot at 0x7fbd02b08160>
```











## Answer

First of all we would use plots in order to visualise the trend in time of each variable alone. Scatter is better to understand the correlation between variables (without time in it).

For data that is empty we should interpolate (RALTs), and not for the rest.

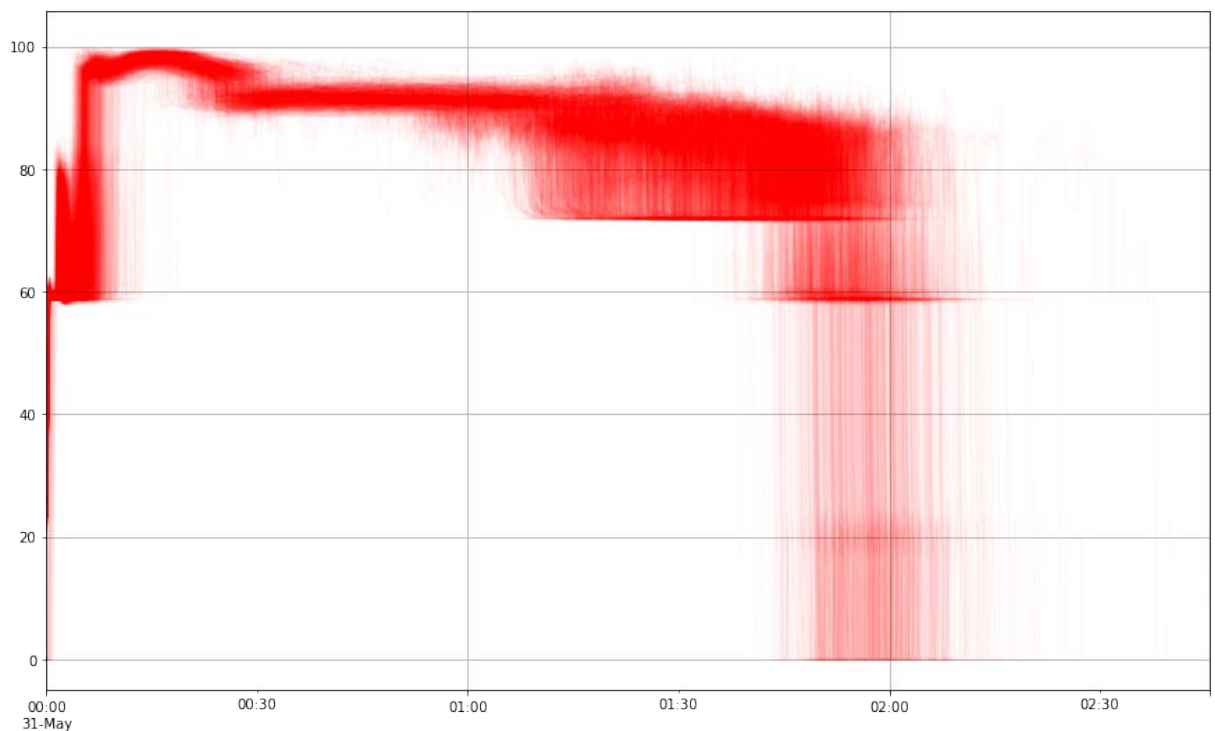
We should interpolate as a solution to the NaN values, but if there are too many NaNs it's not worth it, it may be better to join all three RALTs together

If it is interesting to see the variables for a given flight, it is more informative to view the set of values for all flights in order to understand what are the significant/normal values and what are those which are abnormal.

### Question 2 Visualize N21 variable for all flights

For the `N21` variable, for example, display all of the flights on the same figure. Use `alpha` parameter to add transparency to your plot. Is there any pattern? Comment the variabilities you observe.

```
In [51]: for f in flights:
          f.N21.plot(alpha=0.02, color='red', figsize=[15,9])
```



### Answer

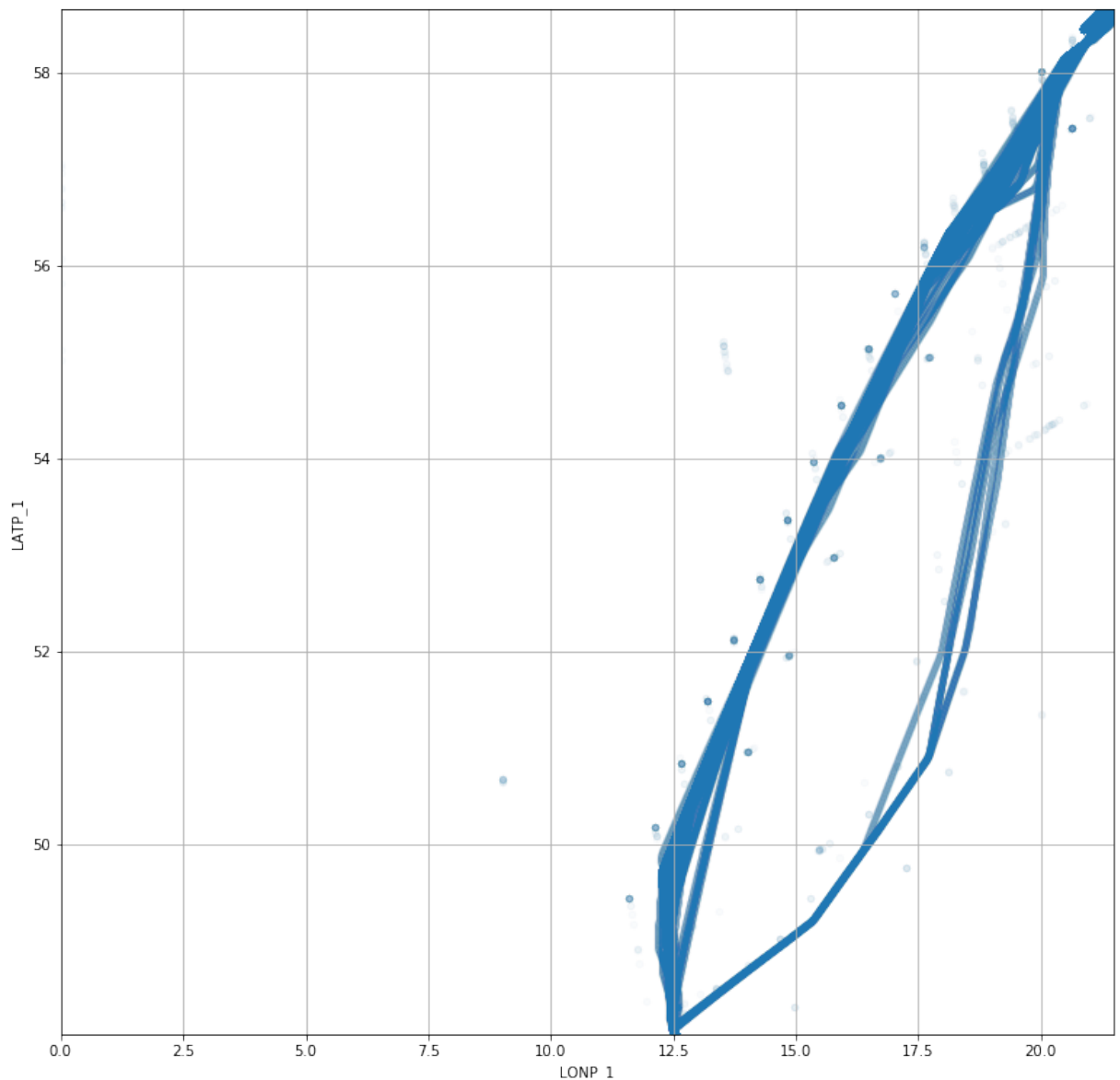
Even though they follow a clear pattern, there is a large variance. Nevertheless, there is a clear trend, with the ascension up to 00:30, cruise from then on and landing at the end.

Some variables must be analyzed together, such as latitude and longitude, otherwise the visualization information will be incomplete, we could be missing something.

### Question 3 Visualize latitude against longitude for all flights

Display the trajectories (LONP\_1, LATP\_1) of a subset of flights, for example 50 flights. What do you see? Keep in mind that the data during the beginning of the recording may be abnormal. What insight do you lose when you plot LONP\_1 against LATP\_1 ?

```
In [49]: # fig, ax = plt.subplots(figsize=[16,16])
# for flight in flights[:50]:
#     ax.plot(x=flight.LONP_1, y=flight.LATP_1, alpha=0.1)
# plt.show()
fig, ax = plt.subplots(figsize=[13, 13])
for f in flights[:50]:
    f[10:].plot(x='LONP_1', y='LATP_1', ax=ax, alpha = 0.02, legend=False,
e, kind='scatter')
    ax.set_xlim(min(f.LONP_1[10:]), max(f.LONP_1[10:]))
    ax.set_ylim(min(f.LATP_1[10:]), max(f.LATP_1[10:])))
```



### Answer

We lose the time, not knowing how fast was the plane travelling. We also lose the height, meaning we have no clue about the takeoff-landing points (although we can have an approximate idea).

Keep in mind that our goal is to understand the nature and the inherent problems of our data, and

its features. Proceed with the visual analysis of the data, looking at different features.

#### Question 4 Recap variables that require pre-processing

Based on your observations as for now, what are the variables requiring processing? For each of these variables, specify the necessary pre-processing required prior to perform data analysis.

#### Answer

- **Longitude & Latitude:** Messy beginning and noisy signals with very high peaks. filtering it out would help.
- **RALTs:** Treatment of NaNs, maybe merging the three measurements.
- **HEAD:** There are some discontinuities, probably due to angle wrapping.
- **NXX and others:** they are quite clean, not need to processing

## Pre-processing

Data pre-processing is essential, in order to separate the errors due to measurement from "normal" data variability, which is representative of the phenomenon that interests us.

#### Question 5 Smooth and filter out abnormal data in trajectories (LATP\_1 and LONP\_1)

Filter the flight trajectories (LATP\_1 and LONP\_1 variables). You can focus on the first 20 flights, that is `flights[:20]`. Display the trajectories before and after smoothing.

```
In [165]: # This is a template code, fill in the blanks, or use your own code

# Give an alias to the first few flights for convenience
fs = flights[:20]

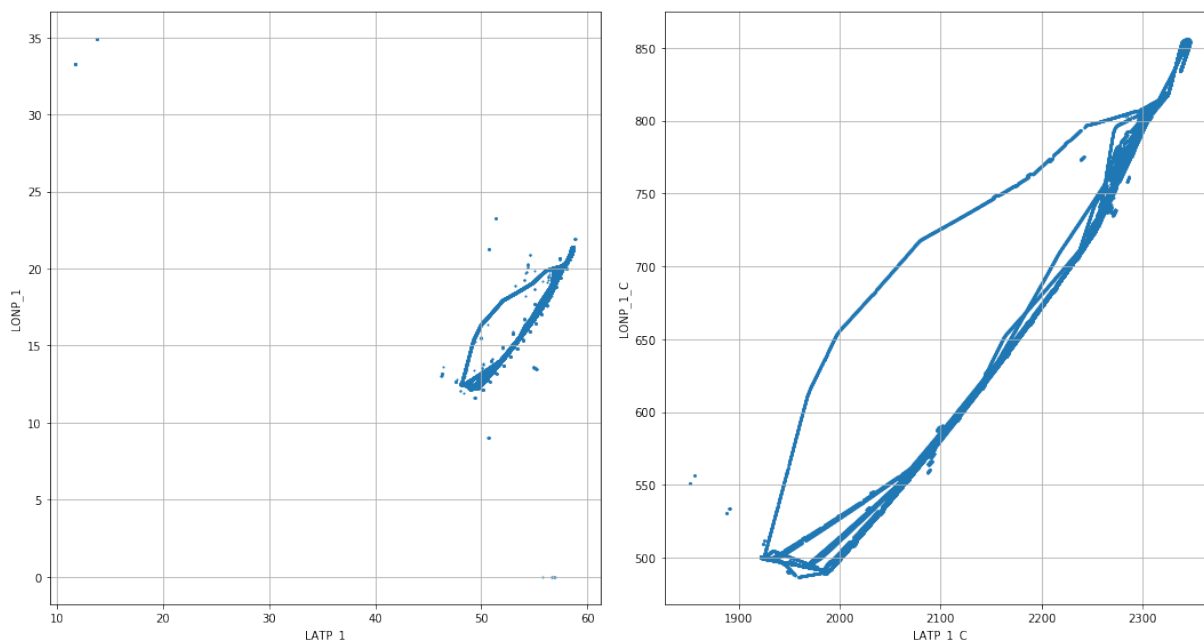
# Set up the figure to plot the trajectories before (ax0) and after smoothing (ax1)
fig, axes = plt.subplots(1, 2, figsize=(15, 8))

# Unpack the axes
ax0, ax1 = axes

# Iterate over fs and add two new smooth columns for each flight
for f in fs:
    f["LATP_1_C"] = f.LATP_1.rolling(window=40).sum() # FILL IN THE BLANKS
    f["LONP_1_C"] = f.LONP_1.rolling(window=40).sum()
```

```
# Iterate over fs and plot the trajectories before and after smoothing
for f in fs:
    # Plot the raw trajectory on ax0
    f.plot(kind="scatter", x="LATP_1", y="LONP_1", s=1, ax=ax0)
    # Plot the smoothed trajectory on ax1
    f.plot(kind="scatter", x="LATP_1_C", y="LONP_1_C", s=1, ax=ax1)

fig.tight_layout()
```



## Answer

Just filtering with a window size of 40 secs and summing does the job

### Question 6 Pre-process HEAD, get rid off discontinuities

Angles are special variables because they "cycle" over their range of values. The `HEAD` variable shows artificial discontinuities: your goal is to eliminate (filter out) such discontinuities. The angle may no longer be between 0 and 360 degrees after the transformation but it will come very handy for some analysis later. Display the data before and after transformation. You can focus on one flight, for example `flights[0]`.

```
In [15]: # Your code goes here ...
```

## Answer

your answer here ...

## Part 2: Analysis

We now turn to the data analysis task. In this part, we will use a **clean** dataset, which has been

prepared for you; nevertheless, the functions you developed in the first part of the notebook can still be used to visualize and inspect the new data. Next, we display the schema of the new dataset you will use:

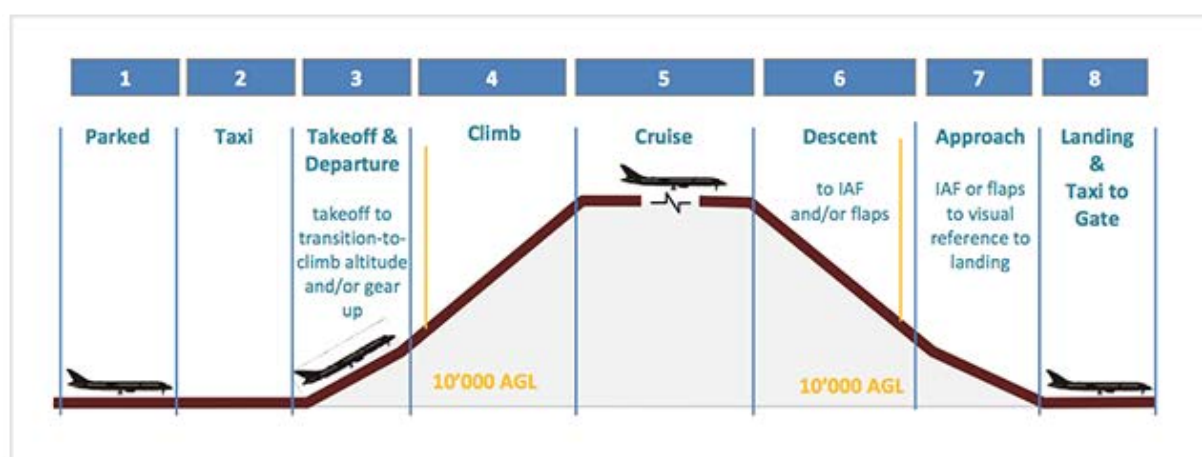
name	description	unit
TIME	elapsed seconds	second
LATP_C	Latitude, Corrected	degree °
LONP_C	Longitude, Corrected	degree °
RALT_F	Radio Altitude, Fused	feet
ALT_STD_C	Relative Altitude, Corrected	feet
HEAD_C	head, Corrected	degree °
HEAD_TRUE	head, without discontinuities	degree °
PITCH_C	pitch, Corrected	degree °
ROLL_C	roll, Corrected	degree °
IAS_C	Indicated Air Speed, Corrected	m/s
N11	speed N1 of the first engine	%
N21	speed N2 of the first engine	%
N12	speed N1 of the second engine	%
N22	speed N2 of the second engine	%
AIR_GROUND	1: ground, 0: air	boolean

Execute the cell below to load the data for part 2

```
In [167]: num_flights = 780
flights = load_data_from_directory(BASE_DIR + "part2/flights/", num_flights)
for f in flights:
    l = len(f)
    new_idx = pd.date_range(start=pd.Timestamp("now").date(), periods=l,
                             freq="S")
    f.set_index(new_idx, inplace=True)
```

There are 780 files in total  
We process 780 files  
100%: [#####]

## Detection of phases of flight



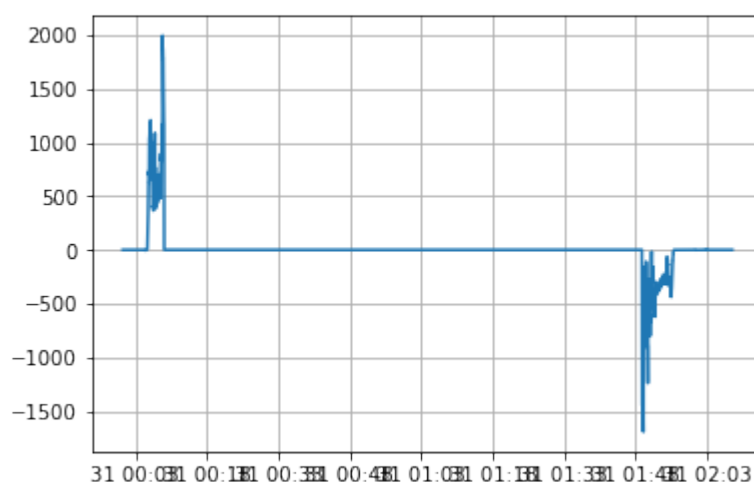
In order to understand the different events that can happen, it is necessary to understand in what phase of the flight the aircraft is located. Indeed, an event that could be regarded as normal in a stage could be abnormal in another stage.

### Question 7 Detect take-off and touch-down phases

Using the clean dataset, detect the take-off phase and the touch-down of all flights. Among all the variables available, what is the variable that tells us the most easily when the take off happens? There is no trap here. Choose the best variable wisely and use it to detect the indices of take-off and touch-down. Plot `ALT_STD_C` 5 mins before and 5 mins after take-off to test your criterion. Do the same for touch-down.

```
In [153]: plt.plot(f.RALT_F.diff(1).rolling(window=20).sum())
```

```
Out[153]: [<matplotlib.lines.Line2D at 0x7fbd05466550>]
```



```
In [154]: fig, ax = plt.subplots(1,2, figsize=[15,9])

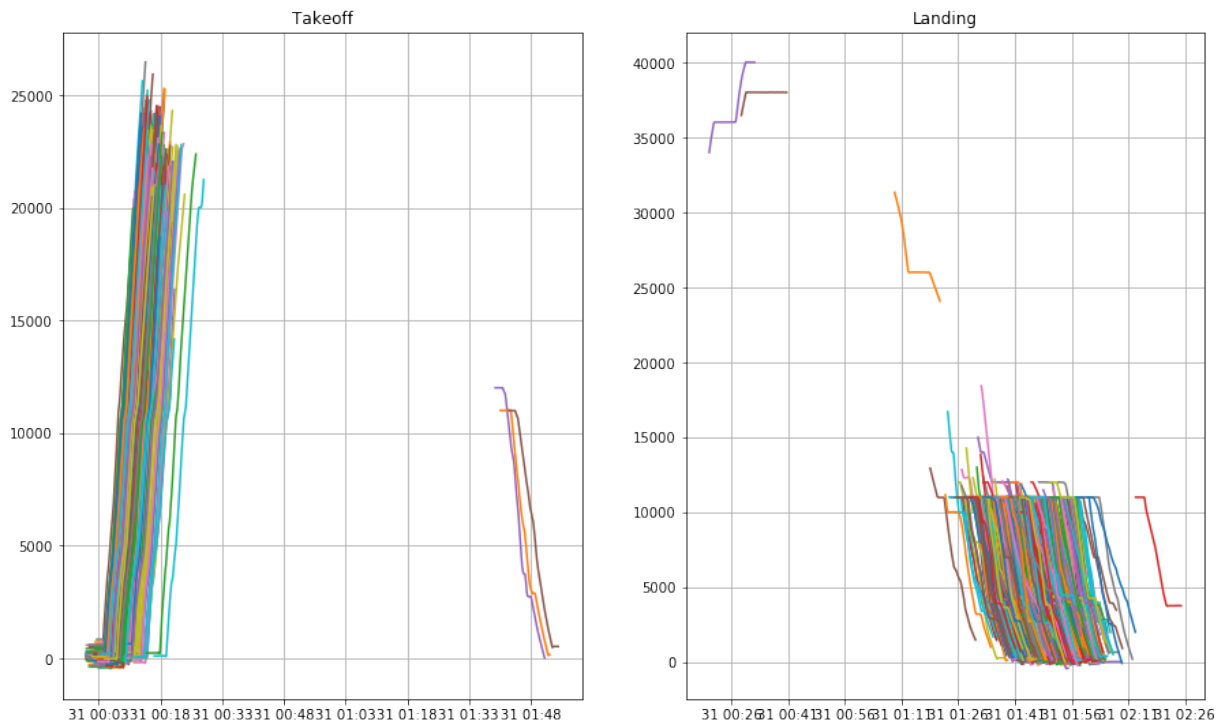
ax[0].set_title('Takeoff')
ax[1].set_title('Landing')

for i,f in enumerate(flights):
    takeoff = f.RALT_F.diff(1).rolling(window=20).sum().idxmax()
```

```

landing = f.RALT_F.diff(1).rolling(window=20).sum().idxmin()
ax[0].plot(f.loc[takeoff-360:takeoff+360].ALT_STD_C)
ax[1].plot(f.loc[landing-360:landing+360].ALT_STD_C)
plt.show()

```



## Answer

Clearly the best variable would have been AIR\_GROUND, a boolean that would tell us the exact moment when it starts flying. However, since we don't have access to it anymore, we are gonna use RALT\_F, and by applying a difference of 5 seconds in the series while detecting the maximum and minimum we obtain the value.

### Question 8 HEAD during take-off and touch-down phases

Plot the `HEAD_C` variable between 20 seconds before the take-off until the take-off itself. Compute the mean of `HEAD_C` during this phase for each individual flight and do a boxplot of the distribution you obtain. Do the same for the touch-down. What do you observe? Is there something significant? Recall [how runways are named] (<https://en.wikipedia.org/wiki/Runway#Naming>)

```

In [155]: fig, ax = plt.subplots(1,2, figsize=[15,9])

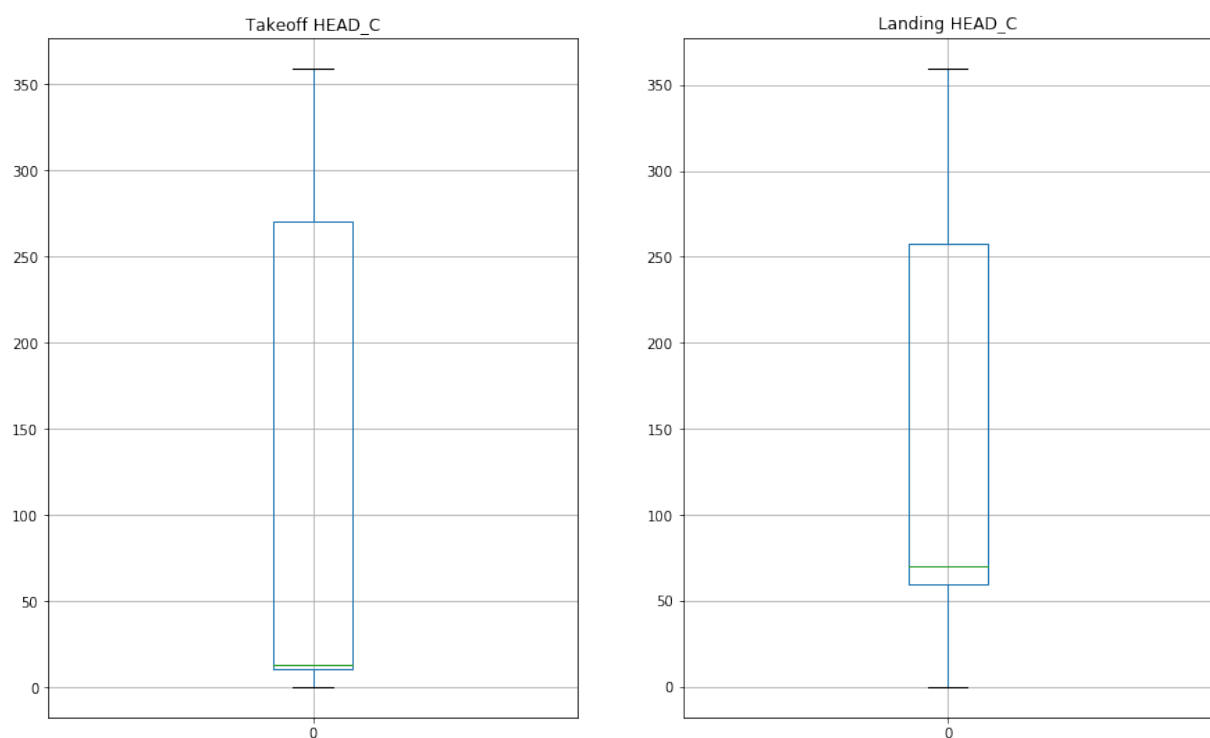
ax[0].set_title('Takeoff HEAD_C')
ax[1].set_title('Landing HEAD_C')
meansTakeoff = []
meansLanding = []
for i,f in enumerate(flights):
    takeoff = f.RALT_F.diff(1).rolling(window=20).sum().idxmax()
    landing = f.RALT_F.diff(1).rolling(window=20).sum().idxmin()
    meansTakeoff.append(f.loc[takeoff-20:takeoff].HEAD_C.mean())
    meansLanding.append(f.loc[landing-20:landing].HEAD_C.mean())
df = pd.DataFrame(meansTakeoff)

```



```
df.plot.box(ax=ax[0])
df2 = pd.DataFrame(meansLanding)
df2.plot.box(ax=ax[1])
```

Out[155]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fbd05e155f8>



## Answer

We see that the vast majority are below 100 degrees. if we divide by 10 we would be getting the most probable names for the runways the plane took.

Also, since the head degree in takeoff is in average smaller than the landing degree, the plane should be turning right in average during his trajectory

Next, we want to detect the moment that the aircraft completed its climb (top of climb) and the moment when the aircraft is in descent phase.

### Question 9 Detect top-of-climb and beginning of descent phases

Plot `ALT_STD_C` a minute before liftoff until five minutes after the top of climb. In another figure plot `ALT_STD_C` a minute before the beginning of descent until the touch-down. For information, a plane is considered:

- in phase of climb if the altitude increases 30 feet/second for 20 seconds
- in stable phase if the altitude does not vary more than 30 feet for 5 minutes
- in phase of descent if the altitude decreases 30 feet/second for 20 seconds

```
In [152]: # This is a template code, fill in the blanks, or use your own code

# Give an alias to flights[0] for convenience
f = flights[0]
```

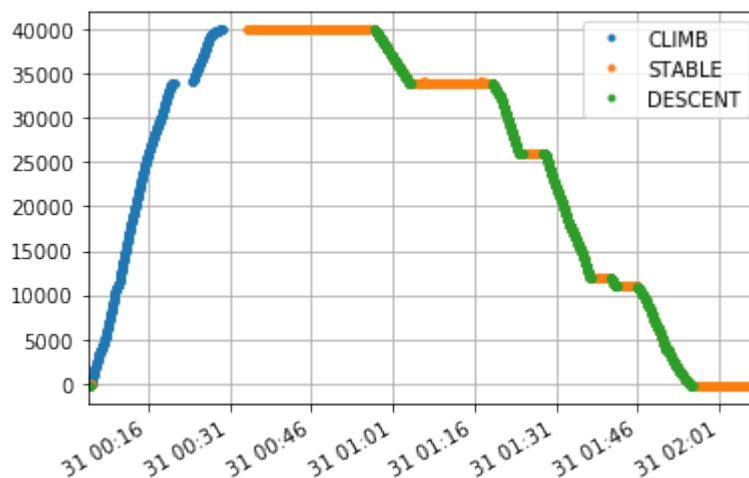
```
f["CLIMB"] = f.ALT_STD_C.diff().rolling(window=20).sum() > 30
f["STABLE"] = f.ALT_STD_C.diff().rolling(window=300).sum() < 30
f["DESCENT"] = f.ALT_STD_C.diff().rolling(window=20).sum() < -30

f[f.CLIMB].ALT_STD_C.plot(color="C0", linestyle="none", marker=".", label="CLIMB") # plot climb phase
f[f.STABLE].ALT_STD_C.plot(color="C1", linestyle="none", marker=".", label="STABLE") # plot stable phase
f[f.DESCENT].ALT_STD_C.plot(color="C2", linestyle="none", marker=".", label="DESCENT") # plot descent phase

top_of_climb = f.CLIMB[-1]
beginning_of_descent = f.DESCENT[0]

plt.legend()
```

Out[152]: <matplotlib.legend.Legend at 0x7fbd067852e8>



**Answer** It works! some data is left out due to window size in stable areas.

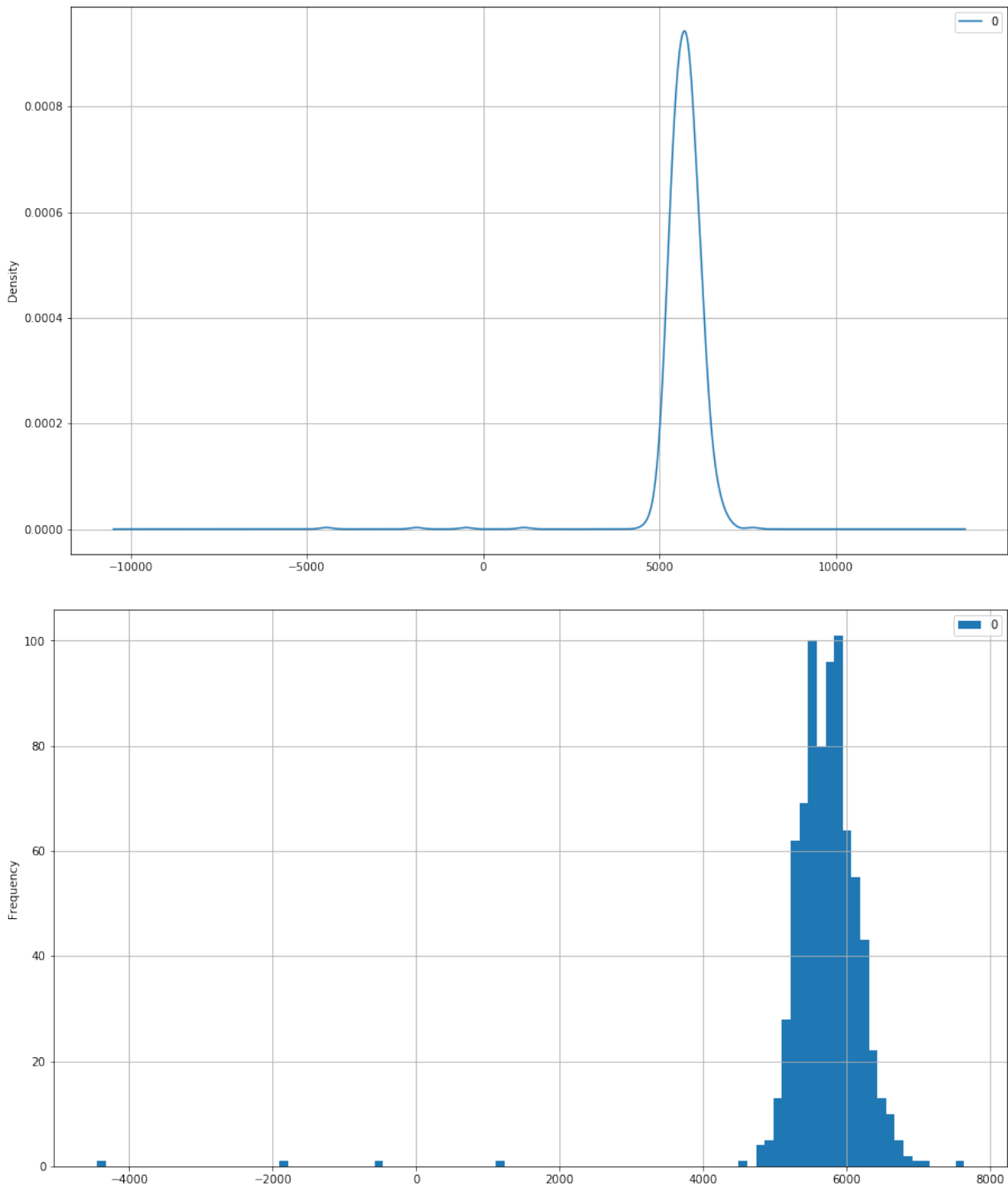
### Question 10 Flight time

Using your criteria to detect the take-off and the touch-down, compute the duration of each flight, and plot the distribution you obtain (boxplot, histogram, kernel density estimation, use your best judgement). Comment the distribution.

```
In [189]: dur_arr = []
for i,f in enumerate(flights):
    takeoff = f.RALT_F.diff(1).rolling(window=20).sum().idxmax()
    landing = f.RALT_F.diff(1).rolling(window=20).sum().idxmin()
    dur_arr.append((landing.value-takeoff.value)/(1e+9))
```

```
In [188]: df = pd.DataFrame(dur_arr)
df.plot.density(figsize=[15,9])
df.plot.hist(figsize=[15,9], bins=100)
```

Out[188]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fbd07d51c50>



**Answer**

Average is arround 6000 secs

**Problems**

Note that the data that we are using in this notebook has been anonymized. This means that the trajectories of a flight have been modified to hide the real information about that flight. In particular, in the dataset we use in this notebook, trajectories have been modified by simple translation and rotation operations

You are asked to find the departure and destination airports of the flights in the dataset. You are guided with sample code to load data from external resources and through several steps that will help you to narrow down the pairs of possible airports that fit with the anonymised data.

We begin by grabbing airport/routes/runways data available on the internet, for example [ourairports](#) (for [airports](#) and [runways](#)) and [openflights](#) (for [routes](#)). These datasets would come useful. You can find the schema of the three datasets below and the code to load the data.

airports.csv

var	description
ident	icao code
type	type
name	airport name
latitude_deg	latitude in °
longitude_deg	longitude in °
elevation_ft	elevation in feet
iata_code	iata code

routes.dat

var	description
AIRLINE	2-letter (IATA) or 3-letter (ICAO) code of the airline.
SOURCE_AIRPORT	3-letter (IATA) code of the source airport.
DESTINATION_AIRPORT	3-letter (IATA) code of the destination airport.

runways.csv

var	description
airport_ident	4-letter (ICAO) code of the airport.
le_ident	low-end runway identity
le_elevation_ft	low-end runway elevation in feet
le_heading_degT	low-end runway heading in °
he_ident	high-end runway identity
he_elevation_ft	high-end runway elevation in feet
he_heading_degT	high-end runway heading in °

The code below has been done for you, it loads the three datasets mentionned and prepare the

pairs dataframe.

```
In [190]: # Load airports data from ourairports.com
airports = pd.read_csv("http://ourairports.com/data/airports.csv",
                        usecols=[1, 2, 3, 4, 5, 6, 13])

# Select large airports
large_airports = airports[(airports.type == "large_airport")]
print("There are " + str(len(large_airports)) +
      " large airports in the world, let's focus on them")

print("airports columns:", airports.columns.values)

# Load routes data from openflights.com
routes = pd.read_csv("https://raw.githubusercontent.com/jpatokal/openflights/master/data/routes.dat",
                      header=0, usecols=[0, 2, 4],
                      names=["AIRLINE", "SOURCE_AIRPORT",
                             "DESTINATION_AIRPORT"])
print("routes columns:", routes.columns.values)

# Load runways data from ourairports.com
runways = pd.read_csv("http://ourairports.com/data/runways.csv", header=
0,
                      usecols=[2, 8, 12, 14, 18],
                      dtype={
                          "le_ident": np.dtype(str),
                          "he_ident": np.dtype(str)
                      })
print("runways columns:", runways.columns.values)

# Create all pairs of large airports
la = large_airports
pairs = pd.merge(la.assign(i=0), la.assign(i=0), how="outer",
                  left_on="i", right_on="i", suffixes=["_origin", "_destination"])

# Compute haversine distance for all pairs of large airports
pairs["haversine_distance"] = pairs.apply(lambda x: haversine((x.latitude_deg_origin,
                                                                x.longitude_deg_origin),
                                                                (x.latitude_deg_destination,
                                                                x.longitude_deg_destination)), axis=1)

del pairs["type_origin"]
del pairs["type_destination"]
del pairs["i"]

pairs = pairs[pairs.ident_origin != pairs.ident_destination]

pairs = pairs.reindex_axis(["ident_origin", "ident_destination", "iata_code_origin", "iata_code_destination",
                            "haversine_distance",
                            "elevation_ft_origin", "elevation_ft_destination",
                            "latitude_deg_origin", "longitude_deg_origin",
                            "latitude_deg_destination", "longitude_deg_destination"], axis=1)

print("pairs columns:", pairs.columns.values)
```

```

There are 577 large airports in the world, let's focus on them
airports columns: ['ident' 'type' 'name' 'latitude_deg' 'longitude_deg'
'elevation_ft'
'iata_code']
routes columns: ['AIRLINE' 'SOURCE_AIRPORT' 'DESTINATION_AIRPORT']
runways columns: ['airport_ident' 'le_ident' 'le_heading_degT' 'he_ident'
' 'he_heading_degT']
pairs columns: ['ident_origin' 'ident_destination' 'iata_code_origin'
'iata_code_destination' 'haversine_distance' 'elevation_ft_origin'
'elevation_ft_destination' 'latitude_deg_origin' 'longitude_deg_origin'
'latitude_deg_destination' 'longitude_deg_destination']

```

Execute the cell below to load the data created by the code above

```

In [191]: airports = pd.read_pickle(BASE_DIR + "part2/airports.pkl")
large_airports = pd.read_pickle(BASE_DIR + "part2/large_airports.pkl")
routes = pd.read_pickle(BASE_DIR + "part2/routes.pkl")
runways = pd.read_pickle(BASE_DIR + "part2/runways.pkl")
pairs = pd.read_pickle(BASE_DIR + "part2/pairs.pkl")

print("There are " + str(len(large_airports)) +
      " large airports in the world, let's focus on them")

# Plot all airports in longitude-latitude plane
plt.scatter(airports["longitude_deg"], airports["latitude_deg"], s=.1)
# Plot large airports in longitude-latitude plane
plt.scatter(large_airports["longitude_deg"], large_airports["latitude_deg"], s=.1)
plt.xlabel("latitude_deg")
plt.ylabel("longitude_deg")
plt.title("All airports (blue) \n large airports (red)")

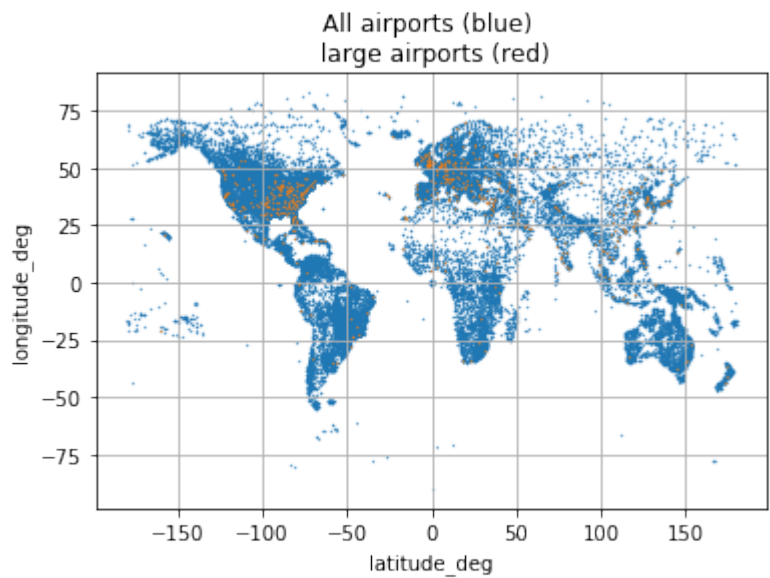
print("airports columns:", airports.columns.values)
print("routes columns:", routes.columns.values)
print("runways columns:", runways.columns.values)
print("pairs columns:", pairs.columns.values)

```

```

There are 574 large airports in the world, let's focus on them
airports columns: ['ident' 'type' 'name' 'latitude_deg' 'longitude_deg'
'elevation_ft'
'iata_code']
routes columns: ['AIRLINE' 'SOURCE_AIRPORT' 'DESTINATION_AIRPORT']
runways columns: ['airport_ident' 'le_ident' 'le_heading_degT' 'he_ident'
' 'he_heading_degT']
pairs columns: ['ident_origin' 'ident_destination' 'iata_code_origin'
'iata_code_destination' 'haversine_distance' 'elevation_ft_origin'
'elevation_ft_destination' 'latitude_deg_origin' 'longitude_deg_origin'
'latitude_deg_destination' 'longitude_deg_destination']

```



You are provided with a dataframe of all pairs of large airports in the world: `pairs`

```
In [38]: pairs.sample(5)
```

Out[38]:

	ident_origin	ident_destination	iata_code_origin	iata_code_destination	haversine
246305	RJTT	EGPF	HND	GLA	9298.6780
293907	VOBL	DNMM	BLR	LOS	8156.3936
73520	KBAD	EGGD	BAD	BRS	7307.9708
58518	FALE	ZBNY	DUR	NAY	11690.466
248634	RKPC	ESMS	CJU	MMX	8266.0264

Question 11.1 Step 1

A first step towards the desanonymisation of the data is to use the distance between the airports. Each entry of ``pairs`` show the latitude and longitude of both airports and the haversine distance between them. Filter the possible pairs of airports by selecting airports that show a distance that is reasonably close to the distance you can compute with the anonymised data. How many pairs of airports do you have left?

```
In [194]: dist_arr = []
for i,f in enumerate(flights):
    takeoff = f.RALT_F.diff(1).rolling(window=20).sum().idxmax()
    landing = f.RALT_F.diff(1).rolling(window=20).sum().idxmin()
    dist_arr.append(haversine((f.loc[takeoff].LATP_C, f.loc[takeoff].LONP_C), (f.loc[landing].LATP_C, f.loc[landing].LONP_C) ))
```

```
In [217]: df= pd.DataFrame(dist_arr)
float(df.mean())
```

```
Out[217]: 1282.3500008648825
```

```
In [226]: pairs2 = pairs[pairs.haversine_distance < float(df.mean())+100]
pairs3 = pairs2[pairs2.haversine_distance > float(df.mean())-100]
pairs3.describe()
```

Out[226]:

	haversine_distance	elevation_ft_origin	elevation_ft_destination	latitude_deg_origi
count	6156.000000	6108.000000	6108.000000	6156.000000
mean	1280.633010	693.090373	693.090373	39.983109
std	58.084858	1036.754591	1036.754591	10.474486
min	1182.391054	-11.000000	-11.000000	-37.673302
25%	1229.778142	80.000000	80.000000	33.663601
50%	1279.448001	353.000000	353.000000	40.467499
75%	1330.138187	823.000000	823.000000	46.238098
max	1382.319408	9205.000000	9205.000000	69.683296

```
In [225]: pairs.describe()
```

Out[225]:

	haversine_distance	elevation_ft_origin	elevation_ft_destination	latitude_deg_origi
count	328902.000000	326037.000000	326037.000000	328902.000000
mean	7390.252178	755.769772	755.769772	33.483940
std	4301.355659	1331.298668	1331.298668	19.342753
min	7.308218	-11.000000	-11.000000	-43.489399
25%	3414.334359	58.000000	58.000000	27.931900
50%	7720.568080	266.000000	266.000000	37.048651
75%	10376.505131	791.000000	791.000000	45.200802
max	19941.928858	10860.000000	10860.000000	69.683296

Answer

We have reduced it to 6K possibilities by limiting it to the mean +-100 in haversine

Question 11.2

Step 2

You should now have a significantly smaller dataframe of possible pairs of airports. The next step is to eliminate the pairs of airports that are connected by commercial flights. You have all the existing commercial routes in the dataset `routes`. Use this dataframe to eliminate the airports that are not connected. How many pairs of airports possible do you have left?



```
In [227]: # This is template code cell, fill in the blanks, or use your own code
selected = pd.merge(pairs3,
                    routes,
                    how='inner',
                    left_on=["iata_code_origin", "iata_code_destination"],
                    right_on=["SOURCE_AIRPORT", "DESTINATION_AIRPORT"])
```

```
In [230]: selected.describe()
```

Out[230]:

	haversine_distance	elevation_ft_origin	elevation_ft_destination	latitude_deg_origi
count	2680.000000	2666.000000	2666.000000	2680.000000
mean	1281.240980	662.489872	670.217929	36.013308
std	57.956149	1178.050412	1208.180925	14.963365
min	1182.656837	-11.000000	-11.000000	-37.673302
25%	1229.681405	36.000000	36.000000	29.993401
50%	1278.940602	202.000000	202.000000	38.133900
75%	1330.667267	723.000000	723.000000	44.882000
max	1382.319408	7841.000000	7841.000000	69.683296

Answer

Further down! up to 2.7K

Question 11.3 Step 3

You now have a list of pairs of airports that are at a reasonable distance with respect to the distance between the airports in the anonymised data and that are connected by a commercial route. We have explored variables in the anonymised data that have not been altered and that may help us to narrow down the possibilities even more. Can you see what variable you may use? What previous question can help you a lot? Choose your criterion and use it to eliminate to pairs of airports that does not fit to the anonymised data.

```
In [ ]: lat_arr1 = []
lat_arr2 = []
for i,f in enumerate(flights):
    takeoff = f.RALT_F.diff(1).rolling(window=20).sum().idxmax()
    landing = f.RALT_F.diff(1).rolling(window=20).sum().idxmin()
    lat_arr1.append(f.loc[takeoff].LATP_C, f.loc[takeoff].LONP_C)
    lat_arr2.append(f.loc[landing].LATP_C, f.loc[landing].LONP_C)
df1 = pd.DataFrame(lat_arr1)
df2 = pd.DataFrame(lat_arr2)
```

```
In [231]: selected[selected.latitude_deg_origin < ...]
```

Out[231]:

	ident_origin	ident_destination	iata_code_origin	iata_code_destination	haversine_distance
0	BIKF	EGPF	KEF	GLA	1347.456794
1	BKPR	EDDF	PRN	FRA	1263.986817
2	BKPR	EDDF	PRN	FRA	1263.986817
3	BKPR	EDDT	PRN	TXL	1251.763447
4	BKPR	LFSB	PRN	BSL	1195.966614

Answer

Clearly using the Latitude and Longitude should help greatly, filtering the same way we did before with the other data. No time to run it!

Question 11.4

Step 4

Is there any other variables that can help discriminate more the airports?

In [15]: `# Your code goes here ...`

Answer

your answer here ...

In [ ]: