

# **Anomaly Detection on Taxi Calls**

We are contacted by a Taxi company:



# **Anomaly Detection on Taxi Calls**

# We are contacted by a Taxi company:

- They have historical data about taxi calls in NYC
- They are interested in detecting "abnormal situations" (so called anomalies)

#### Goals:

- Analyze anomalies (e.g. better size the fleet)
- Anticipate anomalies (so we can prepare)

### Typically referred to as anomaly detection:

- An important industrial problem
- Many context and possible applications

# **Basic Setup**

# Let us start by setting up the notebook:

```
In [1]: %load_ext autoreload
%autoreload 2
#%matplotlib widget
```

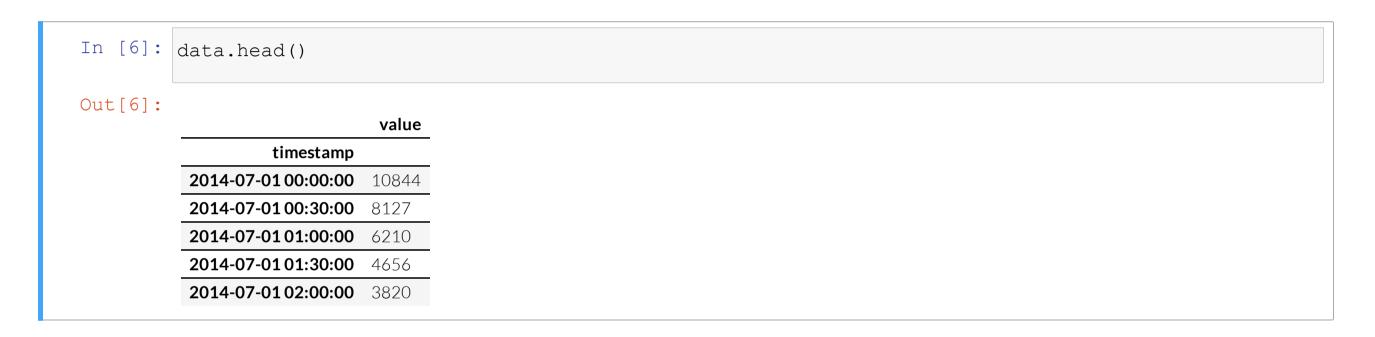
Our module contains a pre-built function to load the data:

```
def load_series(file_name, data_folder):
...
```

- We will use data from the <u>Numenta Anomaly Benchmark (NAB)</u>
- NYC taxi data nyc\_taxi.csv is in the data/realKnownCause folder

```
In [5]: from util import nab # Import our submodule
   data_folder = '../data/nab'
   file_name = 'realKnownCause/nyc_taxi.csv'
   data, labels, windows = nab.load_series(file_name, data_folder)
```

#### Let's have a look at all the data we loaded



data is a pandas DataFrame Object

- It is essentially a table, in this case representing a time series
- There are well defined column names (here "value")
- There is a well defined row index (here "timestamp")
- Jupyter displays DataFrame objects as HTML tables

### **Time Series and Pandas**

#### Our data is a time series

I.e. a sequence whose index represents time

- Specifically, we have a univariate time series...
- ...Since we are tracking only quantity (i.e. one variable)

#### Times series have one difference w.r.t. classical table datasets

- ...l.e. their row index is meaningful
- Since it represents the position of the example in the sequence

# That said, we do not care about how time is represented

- Hence, time series are stored just as usual!
- Their peculiarities arise when we start to manipulate them

### **Time Series and Pandas**

### In pandas:

- Time series are stored as usual, via DataFrame Or Series Objects
- ...You just need to pay more attention to the index

# It may be convenient using a datetime index

- A datetime object in python allows to manipulate dates/hours directly
  - E.g. get year/month/day/hour/minute...
- In pandas they can be used as indices, so that for example:
  - Time stamps are easier to read
  - We can sort rows by time
  - We can represent arbitrarily long gaps between measurements
  - ...

That said, we still deal with normal DataFrame or Series objects

#### Let's have a look at all the data we loaded

Our module contains a function to plot NAB series:

```
In [7]: nab.plot_series(data)
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

■ If are curious, you can look up the <u>function code in the module</u>

#### Let's have a look at all the data we loaded

We can now move to other data structures

### labels is a pandas series object

- Similar to a 1D array
- ...But with a well defined row index

# This series contains the timestamp of all anomalies

■ They are all hand-labeled

#### Let's have a look at all the data we loaded

Let's plot both the series and the labels:

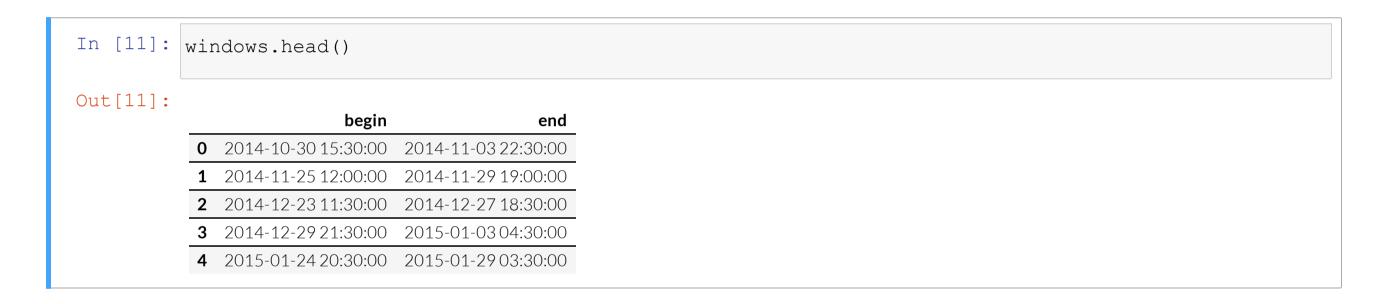
```
In [10]: nab.plot_series(data, labels)
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

Anomalies occur rarely (which is typical for this kind of problem)

#### Let's have a look at all the data we loaded

Now the "windows" data structure:



### windows is a pandas DataFrame Object

- Contains the start/end of windows containing anomalies
- They represent a suitable "resolution" for detecting anomalies
- Reporting the presence of anomalies at any point of the window...
- ...Has some value for the company

#### Let's have a look at all the data we loaded

Let's plot the series, the labels, and the windows all together:

```
In [12]: nab.plot_series(data, labels, windows)
```

A Jupyter widget could not be displayed because the widget state could not be found. This could happen if the kernel storing the widget is no longer available, or if the widget state was not saved in the notebook. You may be able to create the widget by running the appropriate cells.

■ Detections that occur too early/late count as misses