# DATASCIENCE, LEARNING AND APPLICATIONS

## DALAS - EDA

26 janvier 2024

Laure Soulier - Nicolas Baskiotis

# EDA : DEFINITION

## Definitions

### Exploratory Data Analysis (EDA) (Wikipedia)

"Approach of analyzing data sets to summarize their main characteristics, often using statistical graphics and other data visualization methods. A statistical model can be used or not, but primarily EDA is for seeing what the data can tell us beyond the formal modeling and thereby contrasts traditional hypothesis testing."
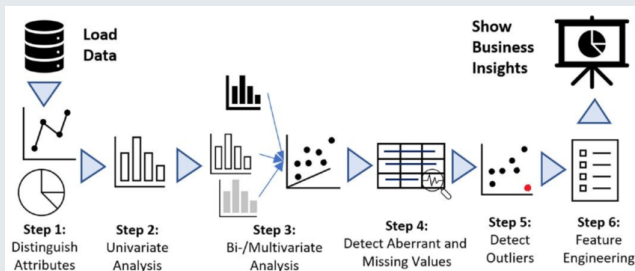


Figure 1 – Source Mahmoud Elansary's thesis - 2021

# EDA vs IDA

SCIENCES
SORBONNE
UNIVERSITÉ

## Exploratory Data Analysis vs. Initial Data Anaysis (Wikipedia)

"EDA is different from **initial data analysis (IDA)**, which focuses more narrowly on checking assumptions required for model fitting and hypothesis testing, and handling missing values and making transformations of variables as needed. EDA encompasses IDA."
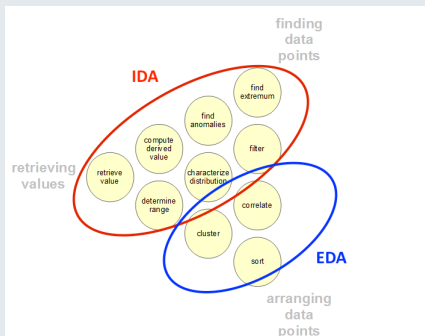


Figure 2 – Source Mahmoud Elansary's thesis - 2021

## Why EDA ? Dirty data : recurrent errors

SCIENCES
SORBONNE
UNIVERSITÉ

- Data might lack coherency or might be poorly collected/done
- Human process : gathering/building/annotating data without objective, without semantics, without consistency



Figure 3 – source https://gallery.dataiku.com/projects/
DKU_CLEANING_CONTACTS/datasets/dss_dirty_data_example/explore/

## Why EDA ? Dirty data : recurrent errors

### Importance of data quality

Impact the representativeness of the model

- Everything under the same label (Name → first name, last name, address → street, apartment, ... )
- No consistency in the same label (zip code is different for the same state)
- No format (phone number, date, ...)
- Spelling mistakes
- Missing values
- Duplication

Example of dataset cleaning : https: //gallery.dataiku.com/projects/DKU_CLEANING_CONTACTS/ recipes/compute_dss_dde_fully_clean/

# Types of data

# Facing all data types : Numerical data

- Different types : continuous (interval, ratio, ...), discrete
- Different temporal dimension : single values, sequential values
- Not the same operations, not the same analysis
- Often stored in Data Frame

# Facing all data types : Visual data (images)

- Matrix of pixels → array of pixels
- 1 pixel : 3 dimensions (x,y,(R,G,B))



{R = 97, G = 134, B=162}

```
1  img= imageio.imread("coco.png")
2  # transformation in 2D, we loose proximity between
       pixels
3  img_2D = img.reshape(-1,3)
```

## Facing all data types : Audio data

- Wave : continuous signal
- Need to be transformed into a series of discrete values.



Figure 4 –
https://huggingface.co/learn/audio-course/chapter1/audio_data

```
1  import matplotlib.pyplot as plt
2  import librosa.display
3
4  plt.figure().set_figwidth(12)
5  librosa.display.waveshow(array, sr=sampling_rate)
```

## Facing all data types : Audio data

- Sampling rate : the number of samples taken in one second, measured in hertz (Hz)
- The **amplitude** of a sound : the sound pressure level at any given timestamp, measured in decibels (dB)



Figure 5 – Frequency spectrum
https://huggingface.co/learn/audio-course/chapter1/audio_data

```
1  # get the amplitude spectrum in decibels
2  amplitude_db = librosa.amplitude_to_db(amplitude)
3  # get the frequency bins
4  frequency = librosa.fft_frequencies(sr=sampling_rate,
     n_fft=len(dft_input))
```

## Facing all data types : Audio data



Figure 6 – Spectrogram
https://huggingface.co/learn/audio-course/chapter1/audio_data

```
1 D = librosa.stft(array)
2 S_db = librosa.amplitude_to_db(np.abs(D), ref=np.max)
3 plt.figure().set_figwidth(12)
4 librosa.display.specshow(S_db, x_axis="time", y_axis="
     hz")
5 plt.colorbar()
```

Facing all data types : Textual data

- Different purposes : categorical data or free text
- Free text can by processed to become a vector - or several (of terms, of topics, of sentiment, ....)

| Nom court | Nom complet | Produit | Arrondissement | Localisation |
|---|---|---|---|---|
| BOBILLOT | MARCHÉ BOBILLOT | Alimentaire | 13 | rue Bobillot, cc |
| PORTE DE VANVES | MARCHÉ AUX PUCES PORTE D... | Puces | 14 | Avenue Marc S |
| PORTE MOLITOR | MARCHÉ PORTE MOLITOR | Alimentaire | 16 | sur le trottoir b |
| CONVENTION | MARCHÉ CONVENTION | Alimentaire | 15 | sur les trottoirs |
| ALESIA | MARCHÉ ALESIA | Alimentaire | 13 | rue de la Glacié |
| LECOURBE | MARCHÉ LECOURBE | Alimentaire | 15 | rue Lecourbe, c |
| COURS DE VINCENNES | MARCHÉ COURS DE VINCENNES | Alimentaire | 12 | terre-plein du c |
| MAUBERT | MARCHÉ MAUBERT | Alimentaire | 5 | place Maubert |
| BELLEVILLE | MARCHÉ BELLEVILLE | Alimentaire | 11 | terre-pleins du |
| TELEGRAPHE | MARCHÉ TELEGRAPHE | Alimentaire | 20 | sur les trottoirs |
| CARRE MARIGNY | MARCHÉ AUX TIMBRES CARRE ... | Timbres | 8 | Angle des aver |
| BEAUVAU - BROC | MARCHÉ BEAUVAU (Brocante) | Brocante | 12 | Place d'Aligre |
| JOURDAN | MARCHÉ JOURDAN | Alimentaire | 14 | bd Jourdan, coi |

## Facing all data types : Structured data (XML, JSON, ...)

- (Semi-)structured dataset
- Dataset loading easy with Pandas :
  - JSON

```
1 data_json=pd.read_json("url")
```

  - XML

```
1 data_xml=request.get("url").content
2 obj=XML2DataFrame(data_xml)
3 xml_dataframe=obj.process_data()
```

  - ...

# Dataset structure

## Describing data structure

- Synthesizing information from a dataset using metrics, tables or graphs
- Describing the dataset structure
  - The size, the type of variables

```
1 data.shape #dimension/size
2 data.info() #size of dataframe, type of data
      by column, used memory
3 data.columns #names of columns
4 data[column_name].dtype #type of column_name
```
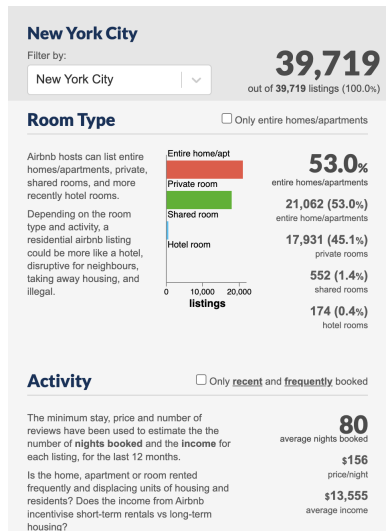
  - Do not hesitate to display some lines in the table

```
1   data.head()
```

## Transforming data

- By default, Pandas uses three main types :
  - integers **int in 32 or 64 bits,**
  - decimal numbers **float in 32 or 64 bits,**
  - Object **objects**, which include most of the other types
- Transforming the type of data
  - Identifying the reason ("price : $48" is an Object and not a float64)
  - Process the data (remove $)
  - Transforming into the right format : `pd.to_numeric()`



**New York City**

Filter by:

New York City

**39,719**
out of **39,719** listings (100.0%)

**Room Type** ☐ Only entire homes/apartments

Airbnb hosts can list entire homes/apartments, private, shared rooms, and more recently hotel rooms.

Depending on the room type and activity, a residential airbnb listing could be more like a hotel, disruptive for neighbours, taking away housing, and illegal.

**53.0%**
entire homes/apartments

**21,062 (53.0%)**
entire home/apartments

**17,931 (45.1%)**
private rooms

**552 (1.4%)**
shared rooms

**174 (0.4%)**
hotel rooms

**Activity** ☐ Only <u>recent</u> and <u>frequently</u> booked

The minimum stay, price and number of reviews have been used to estimate the the number of **nights booked** and the **income** for each listing, for the last 12 months.

Is the home, apartment or room rented frequently and displacing units of housing and residents? Does the income from Airbnb incentivise short-term rentals vs long-term housing?

**80**
average nights booked

**$156**
price/night

**$13,555**
average income

# Fusing/concatenating datasets

SCIENCES
SORBONNE
UNIVERSITÉ

- Join : Aligning two datasets according to a join key, a method (left, right, inner, outer)
- Concatenation : without join key.

## Duplicate data

- Detecting same lines in a dataset and removing duplicate data

```
1 dataset.duplicated().sum() #number of duplicated
    data
2 dataset.duplicated(['NAME']).sum() #focus on the
    column NAME
3 dataset.drop_duplicates(['NAME'],keep="first") #
    remove
```

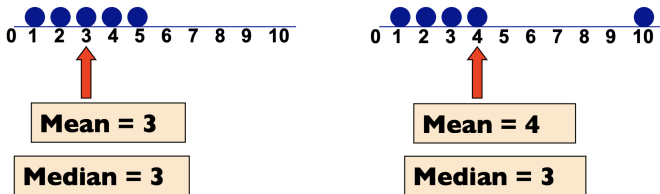| NAME | TITLE | Number |
|---------|--------------|--------|
| Doherty | Officer | 365 |
| Robert | Fire fighter | 457 |
| Robert | Fire Fighter | 127 |
| ... | | |

# DESCRIPTIVE DATA ANALYSIS AND TRANSFORMATION

## Descriptive analysis

- Summary of key characteristics of the data distribution
- Different analyses according to the considered variables :
  - Univariate analysis : on a single variable
  - Bivariate analysis : on two variables
  - Multivariate analysis : on many variables
- Different analyses according to the objectives :
  - Central Tendency measures : general center in which the data are distributed
  - Variability measures : "data spread" or how far away the data are from the center.
  - Relative Standing measures : relative position within the dataset.

## Descriptive statistics on quantitative data

- Mean, Variance, standard deviation, median, percentiles, correlation matrix
  - Mean vs. median : depend on the distribution (outlier/bias, ...)



Figure 7 – (c) Jeffrey Heer - University of Washington

- Based on probability distribution : distribution asymetrie (skewness)

```
1 from scipy.stats import skew
2 skew (dataset["price"])
```

Descriptive analysis on qualitative data

- Modalities, frequency, mode, ...

### Categorial type

This type allows to format data as categories/classes instead of considering just textual data. It allows to **improve the reliability of data** (e.g., modalities are set up and a new data should follow this setup, avoids spelling mistakes in textual data), and **lower the memory consumption**.

```
1 pd.Categorial(["cat1", "cat2", ...])
```

## Visualization

SCIENCES
SORBONNE
UNIVERSITÉ

- Library seaborn https://seaborn.pydata.org/
  - Visualizing relationships : Scatter plot (data order important when lines)

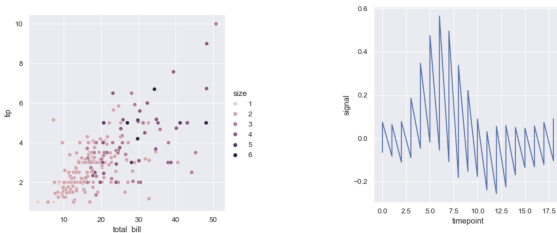

Figure 8 – https://seaborn.pydata.org/tutorial/relational.html

```
1 #relplot function for graphs with Scatter by default
2 sns.relplot(data=tips, x="total_bill", y="tip", hue="
    size", palette="ch:r=-.5,l=.75")
3 #with lines
4 sns.relplot(data=fmri, kind="line", x="timepoint", y="
    signal",estimator=None)
```

## Visualization

- Visualizing distribution
  - Histograms



Figure 9 – https://seaborn.pydata.org/tutorial/distributions.html
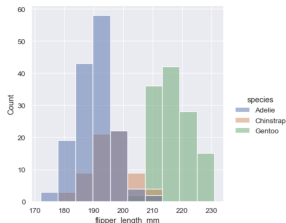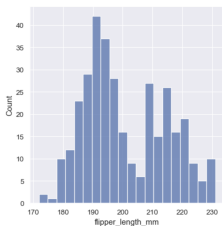
```
1  sns.displot(penguins, x="flipper_length_mm", bins=20)
2  sns.displot(penguins, x="flipper_length_mm", hue="
       species")
```

## Visualization

- Visualizing distribution
  - Kernel density : kernel smooting of probability distribution

$$\hat{f_h}(x) = \frac{1}{n} \sum_{i=1}^{n} K_h(x - x_i) \tag{1}$$



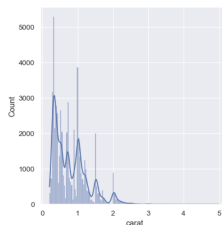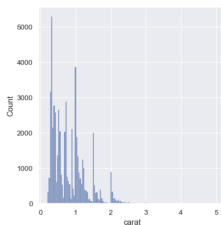Figure 10 – https://seaborn.pydata.org/tutorial/distributions.html

```
1 sns.displot(diamonds, x="carat")
2 sns.displot(diamonds, x="carat", kde=True)
```

## Visualization

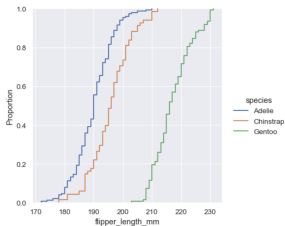- Visualizing distribution
  - Cumulative distribution



Figure 11 – https://seaborn.pydata.org/tutorial/distributions.html

```
1  sns.displot(penguins, x="flipper_length_mm", hue="
       species", kind="ecdf")
```

## Visualization

- Visualizing distribution
    - Bivariate distributions



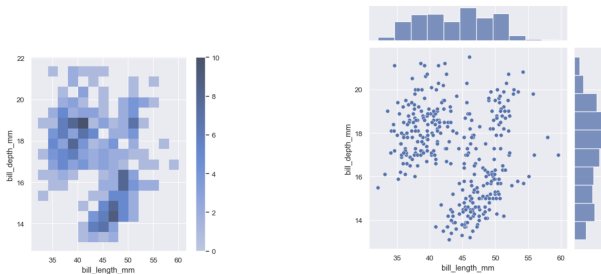Figure 12 – https://seaborn.pydata.org/tutorial/distributions.html

```
1  sns.displot(penguins, x="bill_length_mm", y="
       bill_depth_mm", hue="species")
2  sns.jointplot(data=penguins, x="bill_length_mm", y="
       bill_depth_mm")
```

## Visualization

SCIENCES
SORBONNE
UNIVERSITÉ

- Visualizing distribution
  - Plotting many distribution
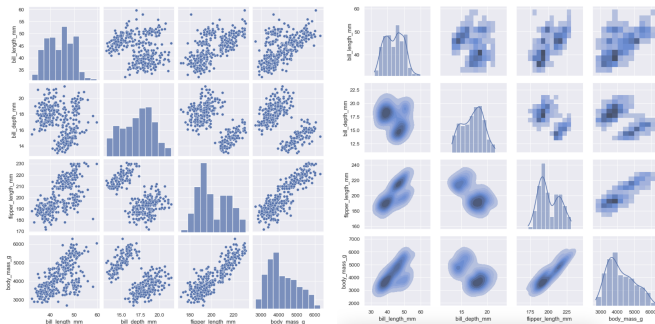


Figure 13 – https://seaborn.pydata.org/tutorial/distributions.html

```
1 sns.pairplot(penguins)
2 g = sns.PairGrid(penguins) #more flexible
3 g.map_upper(sns.histplot)
4 g.map_lower(sns.kdeplot, fill=True)
5 g.map_diag(sns.histplot, kde=True)
```

## Visualization

- Visualizing categorical data
  - Scatter plot



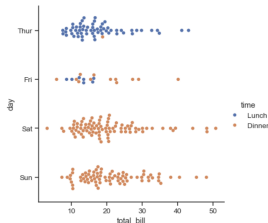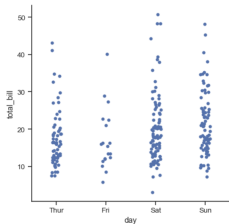Figure 14 – https://seaborn.pydata.org/tutorial/categorical.html

```
1 sns.catplot(data=tips, x="day", y="total_bill")
2 sns.catplot(data=tips, x="total_bill", y="day", hue="
    time", kind="swarm")
```

## Visualization

- Visualizing categorical data
  - Boxplot



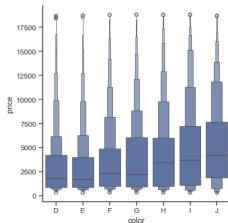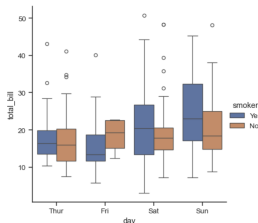Figure 15 – https://seaborn.pydata.org/tutorial/categorical.html

```
1 sns.catplot(data=tips, x="day", y="total_bill", hue="
      smoker", kind="box")
2 sns.catplot(data=diamonds.sort_values("color"), x="
      color", y="price", kind="boxen")
```

## Visualization

- Visualizing categorical data
  - Violin plot



Figure 16 – https://seaborn.pydata.org/tutorial/categorical.html

```
1 sns.catplot(data=tips, x="total_bill", y="day", hue="
      sex", kind="violin")
2 #with data distribution
3 g = sns.catplot(data=tips, x="day", y="total_bill",
      kind="violin", inner=None)
4 sns.swarmplot(data=tips, x="day", y="total_bill",
      color="k", size=3, ax=g.ax)
```

## Visualization

- Visualizing categorical data
  - Violin plot



Figure 17 – https://seaborn.pydata.org/tutorial/categorical.html

```
1  sns.catplot(data=titanic, x="deck", kind="count")
2  sns.catplot(data=titanic, x="sex", y="survived", hue="
       class", kind="bar")
```

## Choose the right plot and the right scale...

## Bivariate analysis

- Correlation (more in 2 weeks)
- Pivot tables : Visualizing/analyzing the intersection of several qualitative variables
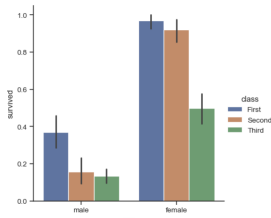
```
1 pd.crosstab(dataset['col1'],dataset['col2'])
```

## Discretization

SCIENCES
SORBONNE
UNIVERSITÉ

- Transforming quantitative variable into a qualitative one.
- Example : age into classes

```
1  pd.cut(dataset["age"],bins=3, labels=range(5)) #
       constant interval
2  pd.cut(dataset["age"],bins=[dataset["age"].min(), 40,
       dataset["age"].max()], include_lowest=True) #
       interval defined by a user
3  pd.qcut(dataset["age"],q-3) # intervals with uniform
       frequency
```

## Missing data

- Identifying why ? (capture, transformation, other ?)
- Deleting observations with missing data
    - Reduce the size of the dataset

```
1 dataset.dropna()
```

- Completing with mean, mode, median for quantitative variables
    - Useful if the missing values occur completely randomly
    - Or in case of rare frequency

```
1 dataset.fillna(dataset[col].mean())
2
3 #autre option avec scikit-learn
4 from sklearn.impute import SimpleImputer
5 imputer=SimpleImputer(strategy="mean")
6 new_dataset=imputer.fit_transform(dataset.
      select_dtypes(np.number))
```

- Add new modality for qualitative variables with .fillna()
- More advances methods (multiple data imputation, KNN)

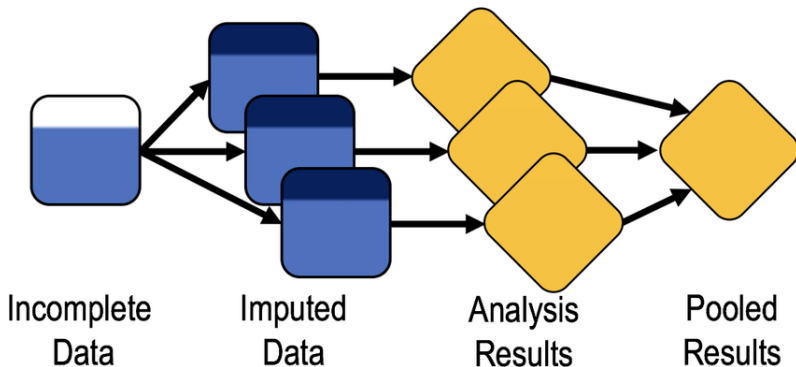## Missing data : multiple imputation



Figure 19 – Nissen et al. 2019

# Outlier detection

- Data that is markedly different from others
- Causes (important to understand why) :
  - Data errors : wrong measurement, wrong annotation, error reporting, ... (1.73 cm for human height, income in billions euros vs. euros)
  - Normal variance in the data : Outside of the 99.7% of the data pointing within three stdev. Those data are legitimate but skew some of the descriptive statistics (e.g., mean).



  - Data from other distribution classes : originate from incorrect assumptions (surge in retail after Thanksgiving vs. daily retail)

## Outlier detection

- Examples
    - Click fraud in online advertising for free internet services
        - Fraudulent traffic does not follow logical actions
        - It contains repetitive actions
        - Signals : Very high click depth, time between each click, high number of clicks in a session, IP different from the target market, ...
    - Credit card fraud
        - Difficult task : irregular purchase is our regular life
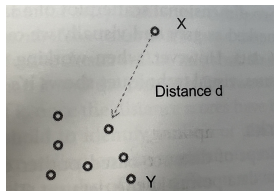        - The recurrence of irregular purchases is a signal

## Anomaly detection techniques

- Statistical methods
    - Normal distribution with parameters estimated on the dataset (mean, stdev). Outliers are detected according on where they fall in the standard normal distribution
- Data mining methods
    - Distance-based : Average distance of the nearest neighbor, outliers will have a higher value than other points



- Clustering : detection with minimum threshold to belong to clusters
- Classification techniques : with dedicated label

# Acceleration and parallelization with Numba and Dask

## Parallelization with Dask

- Dask https://domino.ai/blog/dask-step-by-step-tutorial

```
1 %%time
2 ## Wrapping the function calls using dask.delayed
3 x = delayed(calculate_square)(10)
4 y = delayed(calculate_square)(20)
5 z = delayed(get_sum)(x, y)
6 ## visualize the task graph
7 z.visualize()
```

## Acceleration with Numba

- Python : interpreted language, not optimized
- Parallelization can be adapted to accelerate the code
- If nor sufficient different alternatives :
    - Changing the language (C, C++, Cython : python with a compiler)
    - Use Numba (https://numba.pydata.org/) : does not require to change the python code

```python
1  # Python without Numba : 943 ns + 20.8 ns per loop
2  def hypot_python(x, y) :
3    return math.sqrt(x**2 + y**2)
4
5  # Numba with decorator @jit : 193 ns + 5.56 ns
6  def hypot_numba_jit(x, y) :
7    return math.sqrt(x**2 + y**2)
8
9  # Numba autojit function to transform the Python
      function : 194 ns + 3.56 ns
10 hypot_numba_autojit = autojit(hypot_python)
```

# Before data science…

## Transforming numerical data

- Standard normalization

```
1 from sklearn.preprocessing import StandardScaler
2 scaler=StandardScaler(with_mean=True,with_std=True
      )
3 scaler.fit_transform(dataset)
```

- Change of scale

```
1 from sklearn.preprocessing import MinMaxScaler
2 minmaxScaler=MinMaxScaler((0,100))
3 minmaxScaler.fit_transform(dataset)
```

- Box-Cox transformation : allow to transform data so that it follows Normal law

```
1 from scipy import stats
2 stats.boxcox(dataset["earnings"])
```

## Transforming textual data

1-hot encoding image/explication

```
1 pd.get_dummies(dataset["description"]
2
3 #with scikit-learn
4 from sklearn.preprocessing import OneHotEncoder
5 encode=OneHotEncoder(sparse=False)
6 encode.fit_transform(....)
```

# Sampling

- Random sampling without replacement

```
1 dataset.sample(n=1000)
```

- Stratified sampling

```
1 dataset.groupby('type').apply(lambda x: x.sample(
    frac=.1))
```

## Roadmap for data exploration

- Organize the dataset : structure the dataset with standard rows and columns
- Find the central point for each attribute (mean, mode, ...)
- Understand the spread of attributes (std, range, ...)
- Visualize the distribution of each attribute
- Detect outliers
- Understand the relationships between attributes