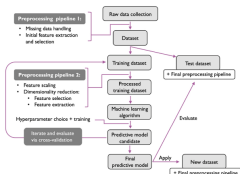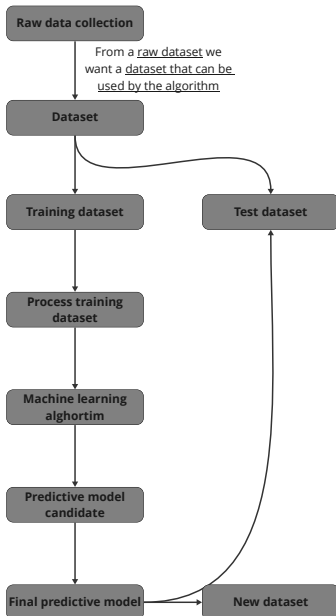MACHINE
LEARNING
WORKFLOW

**Estimator**
- An object which manages the estimation and decoding of a model
- Must provide the methods:
  `fit(data, target)` or
  `fit(data)`
- The model is estimated as a deterministic function of:
  - Data passed to the most recent call to `fit()`, or `partial_fit()`
  - Parameters provided in object construction = hyperparameters or through the method `set_params()`
  - A random state `numpy.random.RandomState`

**Predictor**
- An estimator supporting
  - `predict()`
  - `fit_predict()`

**Transformer**
- An estimator supporting
  - `transform()`: transforms the input into some transformed space. If the transformer was not already fitted, calling this method raise an exception.
  - `fit_transform()`

**Raw data collection**

From a raw dataset we want a dataset that can be used by the algorithm

**Dataset**

**Training dataset**     **Test dataset**

**Process training dataset**

**Machine learning alghortim**

**Predictive model candidate**

**Final predictive model**     **New dataset**

| age | length_snew | hammer_stnen | diameter | rank | size | country |
|---|---|---|---|---|---|---|
| NUMERICAL | NUMERICAL MISSING VALUE | NUMERICAL MISSING VALUE | NUMERICAL | ORDINAL | ORDINAL MISSING VALUE | NOMINAL |
| ⇓ | ⇓ | | ⇓ | ⇓ | ⇓ | ⇓ |
| MIN MAX SCALER | SIMPLE IMPUTER ↓ STANDARD SCALER | | STANDARD SCALER | ORDINAL ENCODER | SIMPLE IMPUTER ↓ ORDINAL ENCODER | ONEHOTENCODER |
| | PIPELINE | | | | PIPELINE | |

# From raw data to a dataset #1

# From raw data to a dataset #2

Two type of data
- Float (number)
- Categorical (text)

Remove all the columns that have less than (ex 6000) items.
*Some columns will still have some missing values*

**Remove all the columns below a treshold**

**Use .simpleImputer()**

Import and use SimpleImputer. Replace missing value with the **mode** (most frequent value) or with the **mean** of the feature column. Useful for categorical data.

Now the dataset is *completely filled*

## 4  Handling categorical data

- Two type of categorical data
- **Ordinal** (have an order. T-shirt size, grade(bad, ok, nice, good), ...) → **Ordinal encoder**
- **Nominal** (no order. Nations, book title...) **OneHotEncoding**

Or skip this if there are no categorical data (very rare)

### 4A. Ordinal data

Repeat for all ordinal data column

**Import OrdinalEncoder**

**Choose the order and the column**

**Double check**

Double check by using playground or playground.io(b) check that the table have new columns with 0 or 1 (one for each country for the column selected

### 4B. Nominal data

Repeat for all nominal data column

**Import OneHotEncoder**

**Start and create temp**

**Create names**

**Merge all**

**Double check**

# From raw data to a dataset #3

```
mm_scaler = MinMaxScaler()
mm_scaler.fit_transform(input_scaled['nb'])[:n]

array([[0.3333],
       [0.45 ],
       [0.5 ],
       ...
       [0.7435],
       [0.5543],
       [0.44 ]])
```

## 5  Handling numerical data

Majority of ML algorithms work better if features are on the same scale. It avoids
- Dominance of some features in the loss computation
- Dominance of some features in Euclidean distance

Most common approaches
- Normalization (or min-max scaling)
- Standardization

## 5A. Normalization or min-max scaling

Apply the min-max scaling to numerical features

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

where $x_{min}$ and $x_{max}$ are the largest and smallest values in the feature vector we are scaling
- Sklearn provides the preprocessing.MinMaxScaler class

```
from sklearn.preprocessing import MinMaxScaler
```

**Import MinMaxScaler**

**Apply MinMaxScaler**

```
scaler = StandardScaler()
scaler.fit_transform(engage[out['input_scaler', 'teacher', 'teacher_strength']])

array([[ 0.5676928,  0,  ],
       [ 1.1354929,  0,  ],
       [ 0.8698848,  0,  ],
       ...
       [-0.8270321,  0,  ],
       [-0.5757312,  0,  ],
       [-0.8395747,  0,  ]])
```

```
from sklearn.preprocessing import StandardScaler
```

**Import StandardScaler**

**Apply StandardScaler**

## 5B. Standardization

More practical for some algorithms with weights initialization
centered around 0

Apply the formula

$$x_{std} = \frac{x - \mu}{\sigma}$$

Where $\mu$ and $\sigma$ are the sample mean and the standard deviation of the feature column
- Sklearn implement standardization by preprocessing.StandardScaler class

## 6. Pipeline

```
my_pipe = Pipeline(
    [ ('input_cc', SimpleImputer(strategy='mean')),
      ('input_scaler', StandardScaler())
    ])
my_pipe.fit_transform(engage[input['input_scaler', 'teacher_strength']])

array([[0,  0,  ],
       [0,  0,  ],
       [0,  0,  ],
       ...
```

```
from sklearn.pipeline import Pipeline
```

**Import Pipeline**

**Define the pipeline**

Choose a name (imp_scaler) and the
composition of the pipeline

# Data loading + confusion matrix #1

## 1. Call the columns

In this way we separate the dataset and the labels

**Load the data and extract the output column**

To get the data that we will use for comparison

**Run column transformer**

## 2. Split the data

We need to split the data between
- Training set
- Test set (used for later)

**Import**
train_test_split

To separate the data

**Use**
train_test_split

train_test_split() uses 4 parameters
- **test_size**: % of observation to put in the test (usually 0.2-0.3)
- **Stratify**: balance between positive and negative instances of a specified column in both train and test size
- **Random state**: shuffle the rows, select portion for test and portion for train

Test_size = train_size = 1

sum(y_train == 1)/len(y_train), sum(y_test == 1)/len(y_test)

Check the average of positive classes in train and test size. **Less the difference better it is**

**Check splitting**

## 2. Confusion matrix



|  | ACTUAL VALUES | |
|---|---|---|
| **PREDICTED VALUES** | **TRUE POSITIVE** | **FALSE POSITIVE** |
| | **FALSE NEGATIVE** | **TRUE NEGATIVE** |

The first diagonal \ express how many items are correctly classified
The second diagonal / express how many items are wrongly classified

(Can be expressed in one row only)

**Fit the perceptron**

**Import**

**Check on train**

**Check on test**

**Accuracy**

More accurate => closer to 1

(less items is precision) (more) (recall) (is more and overfitest, wrong)