

# Adopt a Pet

You are in charge of an animal shelter and you want to predict if the animals you have in your possession can be adopted within 30 days or not.

The dataset at your disposal contains different information about the animals in the shelter: data about the breed or color, data about a cost, data about its health. You even have a short description written by the former owner and a picture of the animal.

We provide you only with the train part and a small test subset so that you can test the whole process.

Deadline: January 15, 2022.

You must submit a zip archive to LMS that contains 3 documents: - A pdf report that outlines the various stages of your work. You will insist on the different hyperparameters of your treatment and for each of them, you will specify on which ranges of values you have tested them. This report will also contain the precision obtained on the train set and on the test set. - the executable notebook containing only the chosen hyper-parameters and not their research. You will leave in this one the execution traces. - A ".joblib" file so that we can execute your code. Of course, the test dataset will be modified and only the predict function of the pipeline will be executed.

The final grade will be based on the quality of the prediction (accuracy score) for 25% and the quality of the work for 75%.

## Table of Contents

- 1 Load train data
  - 1.1 Load the images
  - 1.2 Compute SIFT detector and descriptors for one image
  - 1.3 Extract features and build BOFs
- 2 Build a basic model
- 3 Evaluation of the model

```
In [ ]: import os
        from tqdm import tqdm
```

```
import warnings
warnings.filterwarnings("ignore")

import ssl
ssl._create_default_https_context = ssl._create_unverified_context
```

```
In [ ]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
%matplotlib inline
```

## Load train data

```
In [ ]: path = "https://www.i3s.unice.fr/~riveill/dataset/petfinder-adoption-prediction/"
```

```
In [ ]: breeds = pd.read_csv(path+'breed_labels.csv')
colors = pd.read_csv(path+'color_labels.csv')
states = pd.read_csv(path+'state_labels.csv')

train = pd.read_csv(path+'train.csv')

train['dataset_type'] = 'train'
```

```
In [ ]: len(train)
```

```
Out[ ]: 8168
```

```
In [ ]: # remove to train on the whole set
N = 4000
train = train[:N]
```

```
In [ ]: if 'dataset_type' in train.columns:
```

```
train = train.drop(labels='dataset_type', axis=1)
train.columns
```

```
Out[ ]: Index(['Type', 'Age', 'Gender', 'Color1', 'Color2', 'Color3', 'MaturitySize',
      'FurLength', 'Vaccinated', 'Dewormed', 'Sterilized', 'Health', 'Fee',
      'Description', 'Images', 'Breed', 'target'],
      dtype='object')
```

```
In [ ]: y_train = train['target']
X_train = train.drop(['target'], axis=1)
X_train.head()
```

```
Out[ ]:
```

	Type	Age	Gender	Color1	Color2	Color3	MaturitySize	FurLength	Vaccinated	Dewormed	Sterilized	Health	Fee	Description	Image
0	Cat	12.0	Female	White	Unknown	Unknown	Medium	Yes	Unknown	Unknown	Unknown	Healthy	0.0	We got Luna when she was a kitten in Feb 15'. ...	880e1378' 4.jp
1	Cat	4.0	Male	Golden	White	Unknown	Medium	Yes	No	Yes	No	Healthy	0.0	Ginger Boy was found starving and hungry so I ...	7abe9a0a: 2.jp
2	Cat	12.0	Female	Black	Golden	White	Medium	No	No	No	No	Healthy	0.0	An indoor cat with nice green/ yellowish eyes....	605d07d3: 5.jp
3	Dog	60.0	Male	Black	Gray	White	Medium	No	Yes	Unknown	Unknown	Healthy	0.0	My dog name called boo. He is a male. I feedin...	7ed568ab: 1.jp
4	Cat	36.0	Female	Cream	Gray	White	Large	No	No	No	Yes	Healthy	0.0	1) Foxy is a stray cat which I feed regularly,...	8969b314: 5.jp

```
In [ ]: cat_cols = ['Type', 'Gender', 'Breed', 'Color1', 'Color2', 'Color3',
      'MaturitySize', 'FurLength', 'Vaccinated', 'Dewormed',
```

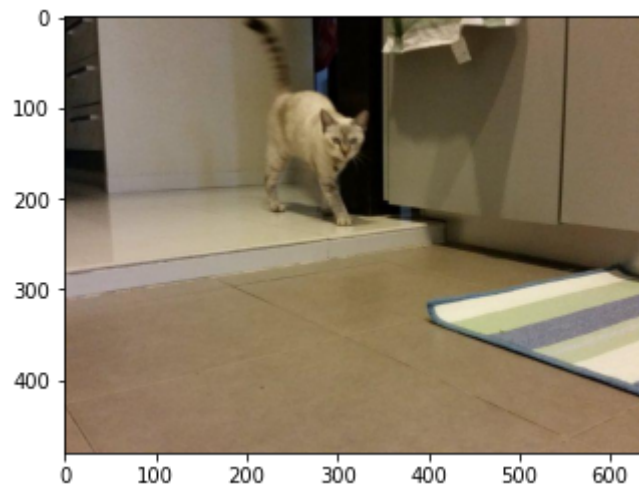
```
'Sterilized', 'Health']  
num_cols = ['Age', 'Fee']  
txt_cols = ['Description']  
img_cols = ['Images']
```

## Load the images

```
In [ ]: # Build the image list of the training set  
img_dir = "train_images/"  
X_train['Images'] = [path+img_dir+img for img in train['Images']]
```

```
In [ ]: from skimage import io  
  
# Read the first image of the list  
img = io.imread(X_train['Images'][0])  
# have a look to the image  
plt.imshow(img)
```

```
Out[ ]: <matplotlib.image.AxesImage at 0x7f1a3d4a9550>
```



## Compute SIFT detector and descriptors for one image

```
In [ ]: # convert the image to grey levels
```

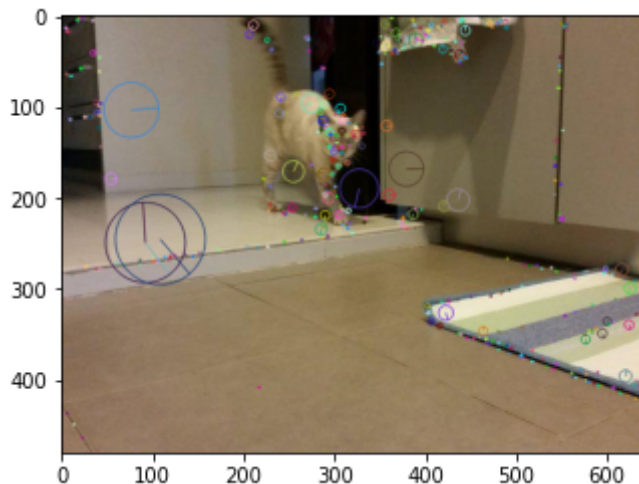
```
import cv2

gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
```

```
In [ ]: # compute SIFT detector and descriptors
sift = cv2.SIFT_create()
kp,des = sift.detectAndCompute(gray,None)

# plot image and descriptors
cv2.drawKeypoints(img,kp,img,flags=cv2.DRAW_MATCHES_FLAGS_DRAW_RICH_KEYPOINTS)
plt.imshow(img)
```

```
Out[ ]: <matplotlib.image.AxesImage at 0x7f1a29ab12e0>
```



## Extract features and build BOFs

```
In [ ]: # First step, extract the SIFTs of each image
# Be carefull: very long process

def extract_SIFT(img_lst):
    nbSIFTs = 0 # Number of SIFTs
    SIFTs = [] # List of SIFTs descriptors
    #dimImgs = [] # Nb of descriptors associated to each images

    for pathImg in tqdm(img_lst, position=0, leave=True):
```

```

img = io.imread(pathImg)
if len(img.shape)==2: # this is a grey level image
    gray = img
else: # we expect the image to be a RGB image or RGBA
    gray = cv2.cvtColor(img, cv2.COLOR_BGR2GRAY)
sift = cv2.SIFT_create()
kp, des = sift.detectAndCompute(gray, None)
if len(kp) == 0 and img.shape[2]==4: #some images are mask on alpha channel: we thus extract this channel if
    gray = img[:, :, 3]
    sift = cv2.SIFT_create()
    kp, des = sift.detectAndCompute(gray, None)

nbSIFTs += des.shape[0]
SIFTs.append(des)
#dimImgs.append(des.shape[0])
return nbSIFTs, SIFTs#, dimImgs

```

```

In [ ]: # nbSIFTs, SIFTs = extract_SIFT(X_train['Images'])
        # print('nbSifts: ', nbSIFTs)

```

```

In [ ]: # Step 2: clusterize the SIFT
        from sklearn.cluster import MiniBatchKMeans

        def clusterize(SIFTs, nb_img_features=5, verbose=False):
            clusterizer = MiniBatchKMeans(n_clusters=nb_img_features) # nb_img_features is a hyperparameter
            # learning of the clustering
            flat_list = SIFTs[0]
            for des in SIFTs[1:]:
                flat_list = np.concatenate((flat_list, des))
                if verbose:
                    print("shape:", des.shape, flat_list.shape)
            clusterizer.fit(flat_list)
            # we now know the label of each SIFT descriptor
            return clusterizer

```

```

In [ ]: # clusterizer = MiniBatchKMeans(n_clusters=5) # nb_img_features is a hyperparameter
        # # learning of the clustering
        # flat_list = SIFTs[0]
        # for des in SIFTs[1:]:

```

```
# flat_list = np.concatenate((flat_list, des))  
# clusterizer.fit(flat_list)
```

```
In [ ]: # clusterizer = clusterize(SIFTs, verbose=True)
```

```
In [ ]: # Step 3: build the BOW representation of each images (i.e. construction of the BOFs)
```

```
def build_BOFs(SIFTs, clusterizer, verbose=False):  
    ok, nok = 0, 0  
    #BOF initialization  
    nb_img_features = clusterizer.get_params()['n_clusters']  
    BOFs = np.empty(shape=(0, nb_img_features), dtype=int)  
  
    # Build label list  
    flat_list = SIFTs[0]  
    for des in SIFTs[1:]:  
        flat_list = np.concatenate((flat_list, des))  
        if verbose:  
            print("shape:", des.shape, flat_list.shape)  
    labels = clusterizer.predict(flat_list)  
  
    # loop on images  
    i = 0 # index for the loop on SIFTs  
    for des in SIFTs:  
        #initialisation of the bof for the current image  
        tmpBof = np.array([0]*nb_img_features)  
        j = 0  
        # for every SIFT of the current image:  
        nbs = des.shape[0]  
        while j < nbs:  
            tmpBof[labels[i]] += 1  
            j+=1  
            i+=1  
        BOFs = np.concatenate((BOFs, tmpBof.reshape(1,-1)), axis=0)  
    if verbose:  
        print("BOFs : ", BOFs)  
  
    return BOFs
```

```
In [ ]: # BOFs = build_BOFs(SIFTs, clusterizer, verbose=True)
```

```
# BOFs.shape
```

```
In [ ]:
```

```
from sklearn.base import BaseEstimator,TransformerMixin

def list_comparaison(l1, l2):
    if not l1 is None \
        and not l2 is None \
        and len(l1)==len(l2) \
        and len(l1)==sum([1 for i,j in zip(l1, l2) if i==j]):
        return True
    return False

class BOF_extractor(BaseEstimator,TransformerMixin):
    X = None
    SIFTs = None
    nbSIFTs = 0

    def __init__(self, nb_img_features=10, verbose=False):
        self.nb_img_features = nb_img_features
        self.verbose = verbose
        self.path = path
        if self.verbose:
            print("BOF.init()")

    def fit(self, X, y=None):
        if self.verbose:
            print("BOF.fit()")
        if list_comparaison(X, self.X):
            SIFTs = self.SIFTs
            nbSIFTs = self.nbSIFTs
        else:
            if self.verbose:
                print("extract_SIFT")
            nbSIFTs, SIFTs = extract_SIFT(X)
        self.X = X
        self.SIFTs = SIFTs
        self.nbSIFTs = nbSIFTs
        self.clusterizer = clusterize(SIFTs, self.nb_img_features, self.verbose)

    def transform(self, X, y=None):
        if self.verbose:
            print("BOF.transform()")
```



```
    if list_comparaison(X, self.X):
        SIFTs = self.SIFTs
        nbSIFTs = self.nbSIFTs
    else:
        if self.verbose:
            print("extract_SIFT")
        nbSIFTs, SIFTs = extract_SIFT(X)

    if self.verbose:
        print("nbSIFTs:", nbSIFTs)
    return build_BOFs(SIFTs, self.clusterizer, self.verbose)

def fit_transform(self, X, y=None):
    if self.verbose:
        print("BOF.fit_transform()")
    if list_comparaison(X, self.X):
        SIFTs = self.SIFTs
        nbSIFTs = self.nbSIFTs
    else:
        if self.verbose:
            print("extract_SIFT")
        nbSIFTs, SIFTs = extract_SIFT(X)
    self.X = X
    self.SIFTs = SIFTs
    self.nbSIFTs = nbSIFTs
    self.clusterizer = clusterize(SIFTs, self.nb_img_features, self.verbose)
    return build_BOFs(SIFTs, self.clusterizer, self.verbose)
```

```
In [ ]: test_BOF_extractor = BOF_extractor(nb_img_features=5, verbose=True)
```

```
BOF.init()
```

```
In [ ]: # test_BOF_extractor.fit(X_train['Images'])
```

```
In [ ]: # BOFs = test_BOF_extractor.transform(X_train['Images'])
        # BOFs.shape
```

```
In [ ]: # BOFs = test_BOF_extractor.fit_transform(X_train['Images'])
        # BOFs.shape
```

```
In [ ]: test = pd.read_csv(path+"test.csv")
y_test = test['target']
X_test = test.drop(['target'], axis=1)

img_dir = "test_images/"
X_test['Images'] = [path+img_dir+img for img in test['Images']]
print(len(X_test))
X_test.head()
```

250

Out[ ]:	Type	Age	Gender	Color1	Color2	Color3	MaturitySize	FurLength	Vaccinated	Dewormed	Sterilized	Health	Fee	Description	
0	Dog	1.0	Male	Black	Gray	Unknown	Medium	Yes	No	No	No	Healthy	0.0	The second puppy of Megan's first litter. Puma...	<a href="https://www">https://www</a>
1	Cat	3.0	Female	Black	Yellow	Unknown	Medium	Yes	No	No	No	Healthy	0.0	Tim is female kitten.active and playful.pls sm...	<a href="https://www">https://www</a>
2	Dog	5.0	Female	Black	Brown	Unknown	Medium	Yes	Yes	Yes	Yes	Healthy	0.0	She was found wearing a red collar, wandering ...	<a href="https://www">https://www</a>
3	Cat	3.0	Male	Cream	Unknown	Unknown	Small	Yes	No	Yes	No	Healthy	0.0	3 months old male kitten. Adopters have to vac...	<a href="https://www">https://www</a>
4	Dog	0.0	Female	Black	Unknown	Unknown	Small	No	Unknown	Unknown	Unknown	Minor Injury	0.0	Please help her. She is an abandon victim. Ver...	<a href="https://www">https://www</a>

```
In [ ]: # BOFs = test_BOF_extractor.transform(X_test['Images'])
        # BOFs.shape
```

## Build a basic model

There are much more interesting things in the dataset and I'm going to explore them, but for now let's build a simple model as a baseline.

```
In [ ]: import os
        from sklearn import set_config

        from sklearn.preprocessing import OneHotEncoder, StandardScaler, MinMaxScaler, FunctionTransformer
        from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
        from sklearn.compose import ColumnTransformer
        from sklearn.pipeline import make_pipeline, FeatureUnion, Pipeline
        from sklearn.decomposition import PCA, SparsePCA, TruncatedSVD

        from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier, AdaBoostClassifier
        from sklearn.metrics import accuracy_score

        set_config(display="text")

        def nb_colors(df):
            """
            Compute the number of known (i.e. not == "Unkown") colors
            """
            return pd.DataFrame((df != "Unknown").sum(axis=1))
```

```
In [ ]: categorical_preprocessor = OneHotEncoder(handle_unknown="ignore", sparse=False)
        numerical_preprocessor = StandardScaler()
        text_preprocessor = TfidfVectorizer()
        image_preprocessor = BOF_extractor(nb_img_features=3, verbose=False)

        nb_colors_transformer = FunctionTransformer(func=nb_colors, validate=False)

        preprocessor = ColumnTransformer(transformers=[
            ("categorical encoding", categorical_preprocessor, cat_cols),
            ("numerical encoding", numerical_preprocessor, num_cols),
            ("text encoding", text_preprocessor, 'Description'),
```

```

    ("image encoding", image_preprocessor, 'Images'),
    ("compute nb colors", nb_colors_transformer,
     [f"Color{i}" for i in range(1, 4)]),
])

dim_red = FeatureUnion([("Truncated SVD", TruncatedSVD(n_components=100))])

preproc_pipe = Pipeline(steps=[
    ("preprocessing", preprocessor),
    ("dimensionality reduction", dim_red),
])
# print("\n--- start preproc_X_train ---")
# preproc_X_train = preproc_pipe.fit_transform(X_train)
# print(f"preproc_X_train shape: {preproc_X_train.shape}")

# print("\n--- start preproc_X_test ---")
# preproc_X_test = preproc_pipe.transform(X_test)
# print(f"preproc_X_test shape: {preproc_X_test.shape}")
# print("\n--- done ---")

--- start preproc_X_train ---
100%|██████████| 4000/4000 [12:46<00:00, 5.22it/s]
preproc_X_train shape: (4000, 100)

--- start preproc_X_test ---
100%|██████████| 250/250 [00:51<00:00, 4.85it/s]
preproc_X_test shape: (250, 100)

--- done ---

```

## Find decent models

In [ ]:

```

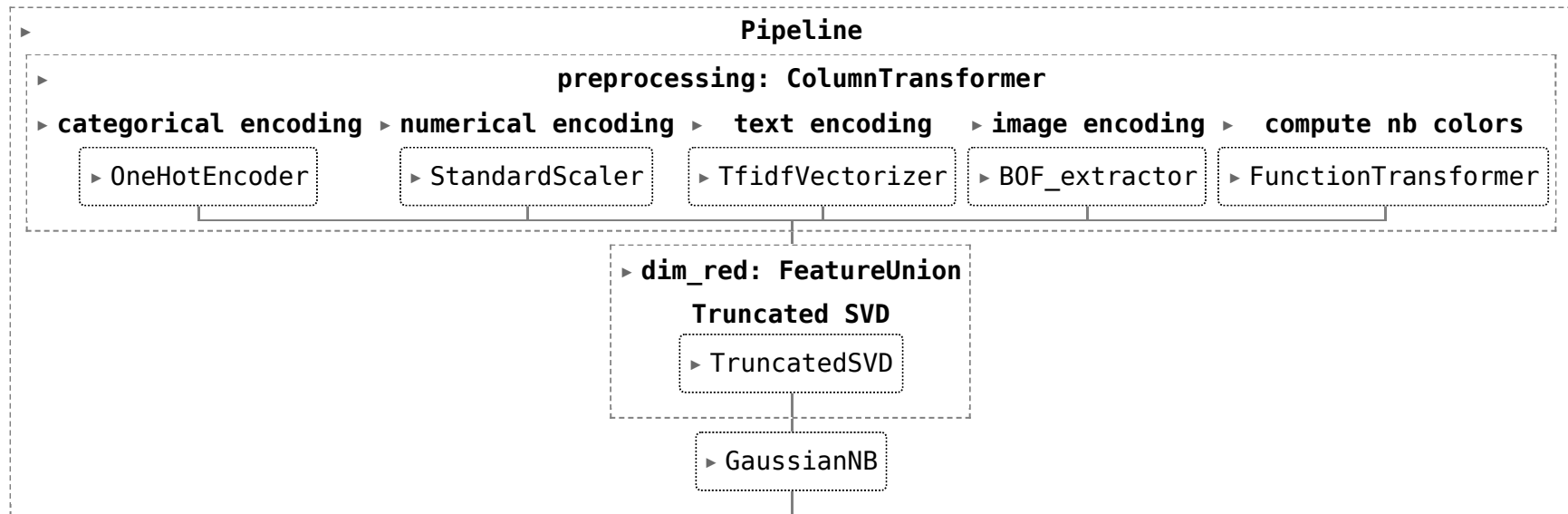
set_config(display="diagram")
set_config(display="text")

model = Pipeline(steps=[
    ("preprocessing", preproc_pipe[0]),
    ("dim_red", preproc_pipe[1]),
    ("classifying", RandomForestClassifier(
        max_depth=200, min_samples_leaf=3, min_samples_split=4)),
])

```

```
model.fit(X_train)
model
```

Out[ ]:



In [ ]:

```
# Save the model
from joblib import dump, load

dump(model, 'limonier.joblib') # Put your name as a model name
```

Out[ ]:

['limonier.joblib']

## Evaluation of the model

We will only execute the following cells.

In [ ]:

```
test = pd.read_csv(path+"test.csv")

y_test = test['target']
X_test = test.drop(['target'], axis=1)

img_dir = "test_images/"
```

```
X_test['Images'] = [path+img_dir+img for img in test['Images']]
print("Test size:", len(X_test))

model = load('limonier.joblib')
y_pred = model.predict(X_train)

print("ACC on train", accuracy_score(y_train, y_pred))
y_pred = model.predict(X_test)
print("ACC on test", accuracy_score(y_test, y_pred))
```

Test size: 250

```
In [ ]: print("ACC on train", accuracy_score(y_train, y_pred))
        y_pred = model.predict(X_test)
        print("ACC on test", accuracy_score(y_test, y_pred))
```

ACC on train 0.58775

100%|██████████| 250/250 [00:50<00:00, 4.97it/s]

ACC on test 0.552

In [ ]:

In [ ]: