

# Revision of Dataframe

## Dataframe from the Previous model

```
In [1]: import pandas as pd
import numpy as np
from sklearn.preprocessing import LabelEncoder
df = pd.read_csv('df.csv')
print("df.shape",df.shape)
item_count = df["target"].value_counts()
print("Number of each species in dataframe:\n",item_count)
```

```
df.shape (3607, 1505)
Number of each species in dataframe:
  A_luchuensis      318
  A_trinidadensis   235
  A_indologenus     232
  A_welwitschiae    223
  A_sclerotiicarbonarius 220
  A_homomorphus     192
  A_ibericus        160
  A_japonicus       141
  A_saccharolyticus 140
  A_niger            140
  A_vadensis        137
  A_costaricaensis  130
  A_heteromorphus   130
  A_carbonarius     120
  A_brasiliensis    120
  A_aculeatinus     120
  A_sclerotioniger  119
  A_ellipticus      110
  A_eucalypticola   110
  A_aculeatus       100
  A_floridensis     90
  A_tubingensis     90
  A_neoniger        80
  A_brunneoviolaceus 80
  A_uvarum          70
Name: target, dtype: int64
```

## Create New Dataframe

```

In [2]: """
Creating individual dataframe of 'A_costaricaensis' , 'A_neoniger' , 'A_tubing
"""

df_class6 = df[df['target'] == 'A_costaricaensis']
df_class16 = df[df['target'] == 'A_neoniger']
df_class22 = df[df['target'] == 'A_tubingensis']
df_class17 = df[df['target'] == 'A_niger']
df_class25 = df[df['target'] == 'A_welwitschiae']

"""

Select 80 data randomly from each dataframe of 'A_costaricaensis' , 'A_neonige
Select 140 data randomly from each dataframe of 'A_niger' and 'A_welwitschiae'
(the number 80,140 based on the minimum of data in each group)
"""

df_class6_rd = df_class6.sample(n = 80)
df_class16_rd = df_class16.sample(n = 80)
df_class22_rd = df_class22.sample(n = 80)
df_class17_rd = df_class17.sample(n = 140)
df_class25_rd = df_class25.sample(n = 140)

"""

Selecting all data from the dataframe in which 'target' is not
'A_costaricaensis','A_neoniger','A_tubingensis', 'A_niger' and 'A_welwitschiae'
"""

target_cut = ['A_costaricaensis','A_neoniger','A_tubingensis','A_niger','A_wel
df_cut = df.loc[~df['target'].isin(target_cut)]
print("df after cut 5 species out.shape",df_cut.shape)

"""

Concatenate dataframe (samples randomly) of
'A_costaricaensis','A_neoniger','A_tubingensis','A_niger','A_welwitschiae'
"""

list_5 = [df_class6_rd,df_class16_rd,df_class22_rd,df_class17_rd,df_class25_rd
df_com5 = pd.concat(list_5, axis=0, ignore_index=True)
print("\nShape of combined 5 species dataframe :",df_com5.shape)
item_counts_5 = df_com5["target"].value_counts()
print("Number of each species :\n",item_counts_5)

"""

Renamed A_costaricaensis and A_neoniger in combined dataframe as A_tubingensis
and renamed A_welwitschiae as A_niger

since A_costaricaensis and A_neoniger are synonyms of A_tubingensis
and A_welwitschiae is a synonyms of A_niger
(Bian et al. 2022)
"""

df_com5_rename = df_com5.replace({'A_costaricaensis':'A_tubingensis',
                                'A_neoniger':'A_tubingensis',
                                'A_welwitschiae':'A_niger'})
print("\nShape of renamed dataframe :",df_com5_rename.shape)
print

item_counts = df_com5_rename["target"].value_counts()

```

```
print("Number of each species :\n",item_counts)

"""
Combine dataframe of cut dataframe and renamed dataframe
"""
list_2 = [df_cut,df_com5_rename]
df_all = pd.concat(list_2, axis=0, ignore_index=True)
df_all.to_csv('dataframe.csv', index=False) # Save to csv file
print("\nSave file as: dataframe.csv")
print("Shape of dataframe :",df_all.shape)

print("Number of each species in dataframe:\n",df_all["target"].value_counts())
```

```
df after cut 5 species out.shape (2944, 1505)
```

```
Shape of combined 5 species dataframe : (520, 1505)
```

```
Number of each species :
```

```
  A_niger          140
```

```
A_welwitschiae    140
```

```
A_costaricaensis   80
```

```
A_neoniger        80
```

```
A_tubingensis      80
```

```
Name: target, dtype: int64
```

```
Shape of renamed dataframe : (520, 1505)
```

```
Number of each species :
```

```
  A_niger          280
```

```
A_tubingensis     240
```

```
Name: target, dtype: int64
```

```
Save file as: dataframe.csv
```

```
Shape of dataframe : (3464, 1505)
```

```
Number of each species in dataframe:
```

```
  A_luchuensis     318
```

```
A_niger            280
```

```
A_tubingensis      240
```

```
A_trinidadensis    235
```

```
A_indologenus      232
```

```
A_sclerotiicarbonarius 220
```

```
A_homomorphus      192
```

```
A_ibericus         160
```

```
A_japonicus        141
```

```
A_saccharolyticus  140
```

```
A_vadensis         137
```

```
A_heteromorphus    130
```

```
A_aculeatinus      120
```

```
A_brasiliensis     120
```

```
A_carbonarius       120
```

```
A_sclerotioniger   119
```

```
A_ellipticus        110
```

```
A_eucalypticola    110
```

```
A_aculeatus         100
```

```
A_floridensis       90
```

```
A_brunneoviolaceus  80
```

```
A_uvarum            70
```

```
Name: target, dtype: int64
```

## Count no. of member in each set & Print class mapping encoder

```
In [3]: cols =[x for x in df_all.columns if x not in ['target']]
rowused = []
for i in range (len(df_all)):
    if i % 10 == 0:
        rowused.append('test')

    elif i % 10 == 1:
        rowused.append('validate')

    else:
        rowused.append('train')

df_all['rowused'] = rowused
dd=df_all['rowused'].sample(len(df_all))
test_set=df_all[df_all['rowused']=='test']
validate_set=df_all[df_all['rowused']=='validate']
train_set=df_all[df_all['rowused']=='train']
print('Count test_set:\n',test_set['target'].value_counts())
print('\nCount validate_set:\n',validate_set['target'].value_counts())
print('\nCount train_set:\n',train_set['target'].value_counts())
print('-----')
label_encoder = LabelEncoder()
data_y = df_all.loc[:, 'target']
encoded_y = label_encoder.fit_transform(data_y.values.ravel())
label_encoder_name_mapping = dict(zip(label_encoder.classes_,label_encoder.transform(label_encoder.classes_)))
print('Mapping of Label Encoded Classes:', label_encoder_name_mapping, sep="\n")
```

```
Count test_set:
  A_luchuensis      32
  A_niger            28
  A_tubingensis     24
  A_trinidadensis   24
  A_indologenus     23
  A_sclerotiicarbonarius 22
  A_homomorphus     20
  A_ibericus        16
  A_saccharolyticus 14
  A_japonicus       14
  A_heteromorphus   13
  A_vadensis        13
  A_aculeatinus     12
  A_brasiliensis    12
  A_carbonarius     12
  A_sclerotioniger  12
  A_ellipticus      11
  A_eucalypticola   11
  A_aculeatus       10
  A_floridensis     9
  A_brunneoviolaceus 8
  A_uvarum          7
Name: target, dtype: int64
```

```
Count validate_set:
  A_luchuensis      32
  A_niger            28
  A_tubingensis     24
  A_trinidadensis   24
  A_indologenus     23
  A_sclerotiicarbonarius 22
  A_homomorphus     19
  A_ibericus        16
  A_saccharolyticus 14
  A_japonicus       14
  A_vadensis        14
  A_heteromorphus   13
  A_aculeatinus     12
  A_brasiliensis    12
  A_carbonarius     12
  A_sclerotioniger  12
  A_ellipticus      11
  A_eucalypticola   11
  A_aculeatus       10
  A_floridensis     9
  A_brunneoviolaceus 8
  A_uvarum          7
Name: target, dtype: int64
```

```
Count train_set:
  A_luchuensis      254
  A_niger            224
  A_tubingensis     192
  A_trinidadensis   187
  A_indologenus     186
  A_sclerotiicarbonarius 176
```

A_homomorphus	153
A_ibericus	128
A_japonicus	113
A_saccharolyticus	112
A_vadensis	110
A_heteromorphus	104
A_aculeatinus	96
A_brasiliensis	96
A_carbonarius	96
A_sclerotioniger	95
A_ellipticus	88
A_eucalypticola	88
A_aculeatus	80
A_floridensis	72
A_brunneoviolaceus	64
A_uvarum	56

Name: target, dtype: int64

-----

Mapping of Label Encoded Classes:

{'A\_aculeatinus': 0, 'A\_aculeatus': 1, 'A\_brasiliensis': 2, 'A\_brunneoviolaceus': 3, 'A\_carbonarius': 4, 'A\_ellipticus': 5, 'A\_eucalypticola': 6, 'A\_floridensis': 7, 'A\_heteromorphus': 8, 'A\_homomorphus': 9, 'A\_ibericus': 10, 'A\_indologenus': 11, 'A\_japonicus': 12, 'A\_luchuensis': 13, 'A\_niger': 14, 'A\_saccharolyticus': 15, 'A\_sclerotiiicarbonarius': 16, 'A\_sclerotioniger': 17, 'A\_trinidadensis': 18, 'A\_tubingensis': 19, 'A\_uvarum': 20, 'A\_vadensis': 21}