

Will there be frogs?

Table of Contents

1. [Introduction](#)
2. [Problem Definition](#)
3. [Data](#)
4. [Problem Approach](#)
5. [Data Management](#)
6. [Model Development](#)
7. [Comparison of Results](#)
8. [Conclusion and Recommendations](#)

Introduction

In many countries, an environmental impact study (EIS) is required to assess the potential impact of actions that significantly affect the quality of the human environment. An EIS is very important but can be an expensive undertaking often involving the deployment of ecological experts to collect data. What if it was possible to do a pre-assessment of a project site? A pre-assessment could give an initial indication of areas to focus on and potentially shorten the the field time required.

Problem Definition

The impact of an infrastructure project on amphibian populations forms part of an EIS.

Can the presence of amphibians species near water reservoirs be predicted using features obtained from GS systems and satellite images?

Data

The data consists of 15 input variables and a multi-label target variable with 7 possible values indicating the presence of a certain type of frog.

Attribute Information

Inputs

- MV categorical
- SR numerical
- NR numerical
- TR categorical
- VR categorical
- SUR1 categorical
- SUR2 categorical
- SUR3 categorical
- UR categorical
- FR categorical
- OR categorical
- BR ordinal
- MR categorical
- CR categorical

Target

- Label 1: The presence of Green frogs
- Label 2: The presence of Brown frogs
- Label 3: The presence of Common toad
- Label 4: The presence of Fire-bellied toad
- Label 5: The presence of Tree frog
- Label 6: The presence of Common newt
- Label 7: The presence of Great Crested newt

Data Source

The dataset was obtained from the UCI Machine Learning Repository: [Amphibians](#)

Marcin Blachnik, Marek SoÅŁtylak, Dominika DÅŁbrowska Predicting presence of amphibian species using features obtained from GIS and satellite images. ISPRS International Journal of Geo-Information 8 (3) pp. 123. MDPI, 2019

Problem Approach

This is a multi-label classification problem. This is not a multi-class classification problem as the classes are not mutually exclusive. Multiple species of amphibians can be found at a site.

Two classification algorithms will be evaluated on dataset and the quality of these models will be assessed. For each algorithm, several models will be developed and the best prediction model will be used to compare the algorithms.

Performance Metrics

The area under the receiver operating characteristic curve (AUC) will be used to evaluate the quality of the models. A perfect classifier will have a ROC AUC equal to 1, whereas a purely random classifier will have a ROC AUC equal to 0.5

Data Management

In [180]:

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from sklearn.model_selection import train_test_split
matplotlib inline
```

The data for both road projects is contained in a single csv file. The repository claims that there are no missing values.

In [181]:

```
amphibians = pd.read_csv('amphibians.csv', sep=';', header=1, index_col='ID')
amphibians.head()
```

Out[181]:

	Motorway	SR	NR	TR	VR	SUR1	SUR2	SUR3	UR	FR	...	BR	MR	CR	Green frogs	Brown frogs	Common toad
ID																	
1	A1	600	1	1	4	6	2	10	0	0	...	0	0	1	0	0	
2	A1	700	1	5	1	10	6	10	3	1	...	1	0	1	0	0	1
3	A1	200	1	5	10	6	10	6	10	4	...	1	0	1	0	0	1
4	A1	300	1	5	0	6	10	2	3	4	...	0	0	1	0	0	0
5	A1	600	2	1	4	10	2	6	0	0	...	5	0	1	0	0	1

5 rows x 22 columns

In [182]:

```
amphibians.dtypes
```

Out[182]:

Motorway	object
SR	int64
NR	int64
TR	int64
VR	int64
SUR1	int64
SUR2	int64
SUR3	int64
UR	int64
FR	int64
OR	int64
BR	int64
MR	int64
CR	int64
Green frogs	int64
Brown frogs	int64
Common toad	int64
Fire-bellied toad	int64
Tree frog	int64
Common newt	int64
Great crested newt	int64
dtype:	object

In [183]:

```
# Check for missing values
amphibians.isna().sum()
```

Out[183]:

Motorway	0
SR	0
NR	0
TR	0
VR	0
SUR1	0
SUR2	0
SUR3	0
UR	0
FR	0
OR	0
BR	0
MR	0
CR	0
Green frogs	0
Brown frogs	0
Common toad	0
Fire-bellied toad	0
Tree frog	0
Common newt	0
Great crested newt	0
dtype:	int64

In [184]:

```
# Determine the number of observations in the dataset
amphibians.shape
```

Out[184]:

(189, 22)

The dataset is quite small with only a 189 observations. A small dataset increases the risk of overfitting.

The dataset will be split into train and test set. The train set will be further split into a train and validation set.

In [185]:

```
# View the number of frogs per species
species = ['Green frogs', 'Brown frogs', 'Common toad', 'Fire-bellied toad', 'Tree frog', 'Common newt', 'Great crested newt']
species_count = amphibians[species].sum().to_frame('count').reset_index()
# Group the values according to the motorway to determine if there are any differences
species_count_motorway = amphibians.groupby('Motorway')[species].apply(lambda x: x.agg(species_count))
```

Out[185]:

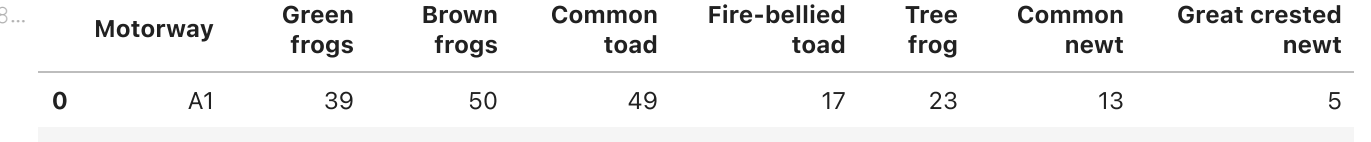
	species	count
6	Great crested newt	21
3	Fire-bellied toad	58
5	Common newt	58
4	Tree frog	71
0	Green frogs	108
2	Common toad	124
1	Brown frogs	148

In [187]:

```
sns.barplot(data=species_count, x='species', y='count')
plt.xticks(rotation=30)
plt.title('The total number of occurrences per species')
```

Out[187]:

Text(0.5, 1.0, 'The total number of occurrences per species')



This figure shows the total number of occurrences per species. Brown frogs were present in 148 of the sites observed, followed by common toads and green frogs. The Great Crested Newt was present at only 21 of the 189 sites observed.

This a form of an unbalanced classification problem because some of the classes are underrepresented.

In [188]:

```
species_count_motorway.reset_index(inplace=True)
species_count_motorway
```

Out[188]:

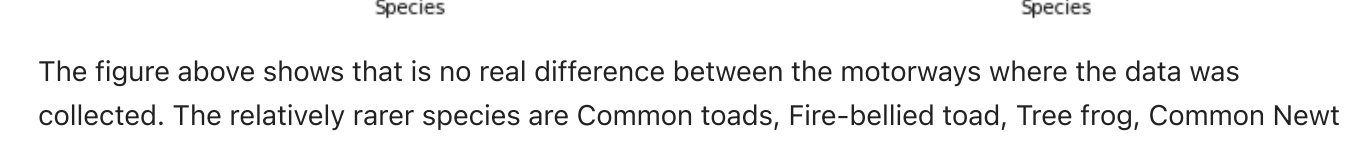
	Motorway	Green frogs	Brown frogs	Common toad	Fire-bellied toad	Tree frog	Common newt	Great crested newt
0	A1	39	50	49	17	23	13	5
1	S52	69	98	75	41	48	45	16

In [189]:

```
g = sns.catplot(data=species_count_motorway, kind='bar', col='Motorway')
g.set_xticklabels(rotation=30)
g.set_axis_labels('Species', 'Count')
```

Out[189]:

<seaborn.axisgrid.FacetGrid at 0x7fe849087c10>



The figure above shows that is no real difference between the motorways where the data was collected. The relatively rarer species are Common toads, Fire-bellied toad, Tree frog, Common Newt and Great crested newts.

In [190]:

```
# Are there any observations with no species present?
frogs_per_observation = amphibians[species].sum(axis=1).to_frame('species_present')
frogs_per_observation[frogs_per_observation['species_present'] == 0]
```

Out[190]:

	species_present
ID	
1	0
6	0
125	0
129	0
132	0
179	0

There are 6 observations where there were no observations any species of amphibians

In [191]:

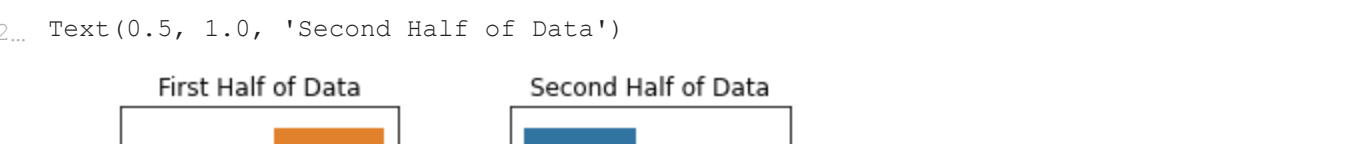
```
# Check whether the data is ordered by splitting the data in half and comparing the 1st
first_half, second_half = train_test_split(amphibians, test_size=0.5, random_state=42)
```

In [192]:

```
fig, ax = plt.subplots(1, 2)
fig.subplots_adjust(hspace=0.4, wspace=0.4)
sns.countplot(ax=ax[0], data=first_half, x='Motorway')
sns.countplot(ax=ax[1], data=second_half, x='Motorway')
ax[0].set_title('First Half of Data')
ax[1].set_title('Second Half of Data')
```

Out[192]:

Text(0.5, 1.0, 'Second Half of Data')



The observations are organized according to the location from where the data was collected (Motorway A1 or S52). We will need to shuffle the data before splitting it so that the model works well for any motorway.

Note: Need to revisit this split to check if the labels are equally distributed between sets

Exploration of the categorical variables

TR - Type of water reservoirs

According to the dictionary, there 10 unique values for this variable

In [193]:

```
amphibians['TR'].value_counts().sort_index()
```

Out[193]:

1	116
2	4
5	12
7	1
11	4
12	23
14	10
15	19
Name: TR, dtype: int64	

In [194]:

```
amphibians.groupby('Motorway')['TR'].value_counts().sort_index()
```

Out[194]:

	Motorway	TR	
	A1	1	48
		2	3
		5	9
		11	4
		12	12
		14	3
		15	1
	S52	1	68
		5	3
		7	1
		12	11
		14	7
		15	18
Name: TR, dtype: int64			

SUR1 - Surroundings 1

According to the dictionary, there 9 unique values for this variable

In [195]:

```
amphibians['SUR1'].value_counts().sort_index()
```

Out[195]:

1	43
2	70
4	5
6	19
7	20
9	5
10	30
14	1
Name: SUR1, dtype: int64	

There are levels of this category that are missing.

SUR2 - Surroundings 2¶

According to the dictionary, there 9 unique values for this variable

In [196]:

```
amphibians['SUR2'].value_counts().sort_index()
```

Out[196]:

1	36
2	41
6	39
7	18
9	10
10	44
11	1
Name: SUR2, dtype: int64	

There are levels of this category that are missing.

SUR3 - Surroundings 3¶

According to the dictionary, there 9 unique values for this variable

In [197]:

```
amphibians['SUR3'].value_counts().sort_index()
```

Out[197]:

1	29
2	29
5	2
6	55
7	18
9	10
10	45
11	1
Name: SUR3, dtype: int64	

CR - Type of shore

In [198]:

```
amphibians['CR'].value_counts().sort_index()
```

Out[198]:

1	186
2	3
Name: CR, dtype: int64	

There are only 3 observations where the shore is concrete. Could it be strong feature for there won't be any amphibians? Or not informative at all?

VR - Intensity of vegetation development

According to the dictionary, there should be 5 unique values for this variable

In [199]:

```
amphibians['VR'].value_counts().sort_index()
```

Out[199]:

0	30
1	25
2	35
3	41
4	28
Name: VR, dtype: int64	

All the levels of the variable appear in the dataset

MR - Maintenance status of the reservoir

According to the dictionary, there should be 3 unique values for this variable

Comment: Trash caused devastation of the reservoir ecosystem. Backfilling and leveling of water reservoirs with ground and debris should also be considered

In [200]:

```
amphibians['MR'].value_counts().sort_index()
```

Out[200]:

0	184
1	4
2	1
Name: MR, dtype: int64	

My gut says this feature won't be informative.

FR - The presence of fishing

According to the dictionary, there should be 3 unique values for this variable

In [201]:

```
amphibians['FR'].value_counts().sort_index()
```

Out[201]:

0	125
1	16
2	15
3	18
4	15
Name: FR, dtype: int64	

There 5 levels of this variable in the dictionary. Check if data was collected differently for the two roads?

In [202]:

```
amphibians['BR'].value_counts().sort_index()
```

Out[202]:

0	39
1	62
2	29
5	46
9	7
10	6
Name: BR, dtype: int64	

The BR or Building development variable indicated the minimum distance from a water reservoir to roads. The options according to the data dictionary are:

1. less than 50 m
2. 50-100 m
3. 100-200 m
4. 200-500 m
5. 500-1000 m
6. >1000 m

In [203]:

```
amphibians['RR'].value_counts().sort_index()
```

Out[203]:

0	47
1	50
2	39
5	42
9	7
10	4
Name: RR, dtype: int64	

The RR or Road variable indicated the minimum distance from a water reservoir to roads. The options according to the data dictionary are:

1. less than 50 m
2. 50-100 m
3. 100-200 m
4. 200-500 m
5. 500-1000 m
6. >1000 m

The values of the two ordinal variables does not correspond to the descriptions of the variables but there are also 6 unique values.

OR - Access from water table to land habitats

According to the dictionary, there should be 4 unique values for this variable

In [204]:

```
amphibians['OR'].value_counts().sort_index()
```

Out[204]:

25	7
50	16
75	22
80	1
99	2
100	141
Name: OR, dtype: int64	

In [205]:

```
# Drop the Motorway column from the dataset as this is not used
amphibians = amphibians.drop(columns='Motorway')
```

In [206]:

```
# Split the data into features and target
amphibians_feat = amphibians.drop(species, axis=1)
```

In [207]:

```
cat_feat_list = ['TR', 'SUR1', 'SUR2', 'SUR3', 'CR', 'VR', 'MR', 'UR', 'FR', 'BR', 'RR']
```

Out[207]:

```
amphibians_feat[cat_feat_list] = amphibians_feat[cat_feat_list].astype('category')
```

In [209]:

```
amphibians_feat.dtypes
```

Out[209]:

SR	int64
NR	int64
TR	category
VR	category
SUR1	category
SUR2	category
SUR3	category
UR	category
FR	category
OR	category
BR	category
MR	category
CR	category
dtype:	object

Model Development

In [210]:

```
from sklearn.preprocessing import OneHotEncoder
from sklearn.pipeline import make_pipeline
from sklearn.pipeline import Pipeline
from sklearn.compose import make_column_transformer
from sklearn.compose import make_column_selector
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import roc_auc_score, roc_curve
```

Overfitting strategy

The original dataset is split into a train and test set.

In [211]:

```
# Shuffle and split the data into train, test and validation sets
X_train, X_test, y_train, y_test = train_test_split(amphibians_feat, amphibians[species])
```

Preprocessing

The categorical columns are converted using one hot encoding

In [212]:

```
one_hot_encoder = make_column_transformer((OneHotEncoder(sparse=False, handle_unknown='raise',
                                                         remainder='passthrough'),
                                                         aucs_scores_macro_rf)
```

Random Forest (RF) Classifier

In [213]:

```
rf = Pipeline(steps=[('one_hot_encoder', one_hot_encoder), ('classifier', RandomForestClassifier())])
```

RF Hyperparameter Tuning

In [214]:

```
param_grid = {'classifier__n_estimators': [5, 10, 15, 20, 30, 40],
              'classifier__max_features': ['auto', 'sqrt', 'log2'],
              'classifier__max_depth': [4, 5, 6, 7, 8],
              'classifier__min_samples_split': [2, 5, 10, 15, 20, 30, 40]}
```

In [215]:

```
gridsearch_rf = GridSearchCV(rf, param_grid, cv=5, scoring='roc_auc', return_train_score=False)
```

In [216]:

```
gridsearch_rf.fit(X_train, y_train)
```

Out[216]:

```
GridSearchCV(cv=5,
             estimator=Pipeline(steps=[('one_hot_encoder',
                                       OneHotEncoder(handle_unknown='ignore', sparse=False),
                                       ColumnTransformer(remainder='passthrough', transformers=[('onehotencoder', OneHotEncoder(handle_unknown='ignore', sparse=False), aucs_scores_macro_rf)]))],
             n_jobs=-1,
             pre_dispatch='all_threads',
             refit=True,
             return_train_score=False,
             scoring='roc_auc',
             verbose=0)
```

Out[217]:</