

BiSSe - Binary State Speciation and Extinction

The model described is a Birth-Death model with two interacting species and the possibility of transitions between them. This model captures the dynamics of two species populations over time, incorporating birth, death, and transition rates.

Parameters

- λ_1, λ_2 : Birth rates of species 1 and species 2, respectively.
- μ_1, μ_2 : Death rates of species 1 and species 2, respectively.
- p_{12}, p_{21} : Transition rates from species 1 to species 2 and from species 2 to species 1, respectively.
- ini_1, ini_2 : Initial populations of species 1 and species 2, respectively.
- T : Maximum time for the simulation.

Model Dynamics

1. **Initialization:** The initial populations of species 1 and species 2 are set based on the given parameters `ini1` and `ini2`. The current time is initialized to zero.
2. **Event Simulation:** The process continues in a loop until the current time exceeds the maximum time `T` or both species' populations become zero.

- **Total Rate Calculation:** At each step, the total rate of events is calculated as the sum of all possible events' rates:

$$\text{total_rate} = n_1(\lambda_1 + \mu_1 + p_{12}) + n_2(\lambda_2 + \mu_2 + p_{21})$$

where n_1 and n_2 are the current populations of species 1 and species 2, respectively.

- **Event Time Sampling:** The time until the next event is sampled from an exponential distribution with the rate parameter `total_rate`.

- **Event Type Sampling:** The type of event is determined by sampling from a discrete distribution with probabilities proportional to the rates of each event:

$$\text{event_probs} = \left[\frac{n_1 \lambda_1}{\text{total_rate}}, \frac{n_2 \lambda_2}{\text{total_rate}}, \frac{n_1 \mu_1}{\text{total_rate}}, \frac{n_2 \mu_2}{\text{total_rate}}, \frac{n_1 p_{12}}{\text{total_rate}}, \frac{n_2 p_{21}}{\text{total_rate}} \right]$$

- **Event Execution:** Based on the sampled event type, the populations are updated accordingly:
 - **Birth of species 1:** $(n_1 \leftarrow n_1 + 1)$
 - **Birth of species 2:** $(n_2 \leftarrow n_2 + 1)$
 - **Death of species 1:** $(n_1 \leftarrow n_1 - 1)$
 - **Death of species 2:** $(n_2 \leftarrow n_2 - 1)$
 - **Transition from species 1 to species 2:** $(n_1 \leftarrow n_1 - 1), (n_2 \leftarrow n_2 + 1)$
 - **Transition from species 2 to species 1:** $(n_2 \leftarrow n_2 - 1), (n_1 \leftarrow n_1 + 1)$
- **Event Recording:** Each event, along with the current time and updated populations, is recorded.

3. **Termination:** The process stops when the current time exceeds the maximum time **T** or both species' populations reach zero.

Output

The function returns a list of events, each represented as a tuple **(time, n1, n2)**, where **time** is the time of the event, and **n1** and **n2** are the populations of species 1 and species 2 after the event.

$$\text{total_rate} = n_1(\lambda_1 + \mu_1 + p_{12}) + n_2(\lambda_2 + \mu_2 + p_{21})$$

```
In [ ]: import numpy as np

def bisse(lam1, lam2, mu1, mu2, p12, p21, ini1, ini2, T, limit_event_size = 1000):
```

```
n1 = ini1.copy()
n2 = ini2.copy()
current_time = 0
events = []
events_list = np.array([1,2,3,4,5,6])
final_T = T

while current_time < T:

    total_population = n1 + n2
    if total_population == 0:
        break
    if len(events) > limit_event_size:
        final_T = current_time
        break

    total_rate = n1*(lam1+mu1+p12) + n2*(lam2+mu2+p21)
    sampled_time = np.random.exponential(1/total_rate)
    current_time += sampled_time

    if current_time > T:
        break

    event_probs = np.array([n1*lam1, n2*lam2, n1*mu1, n2*mu2, n1*p12, n2*p21])/total_rate
    event = np.random.choice(events_list, p=event_probs)

    match event:
        case 1: # specie 1 gives birth
            n1 += 1
        case 2: # specie 2 gives birth
            n2 += 1
        case 3: # specie 1 dies
            n1 -= 1
        case 4: # specie 2 dies
            n2 -= 1
```

```

        case 5: # specie 1 transitions to specie 2
            n1 -= 1
            n2 += 1
        case 6: # specie 2 transitions to specie 1
            n2 -= 1
            n1 += 1
        case _:
            raise ValueError("Invalid event")

    events.append((current_time, n1, n2))

    return final_T, events

```

We randomize the parameters and make an experiment by computing BiSSE and then plotting the evolution of the species over time

```

In [ ]: # parameters
lam1, lam2 = np.random.uniform(0, 1, size=2)
mu1, mu2 = np.random.uniform(0, 0.8, size=2)
p12, p21 = np.random.uniform(0, 0.5, size=2)
max_time = 10
max_num_initial_population = 5
ini1, ini2 = np.random.randint(0, max_num_initial_population, size=2)

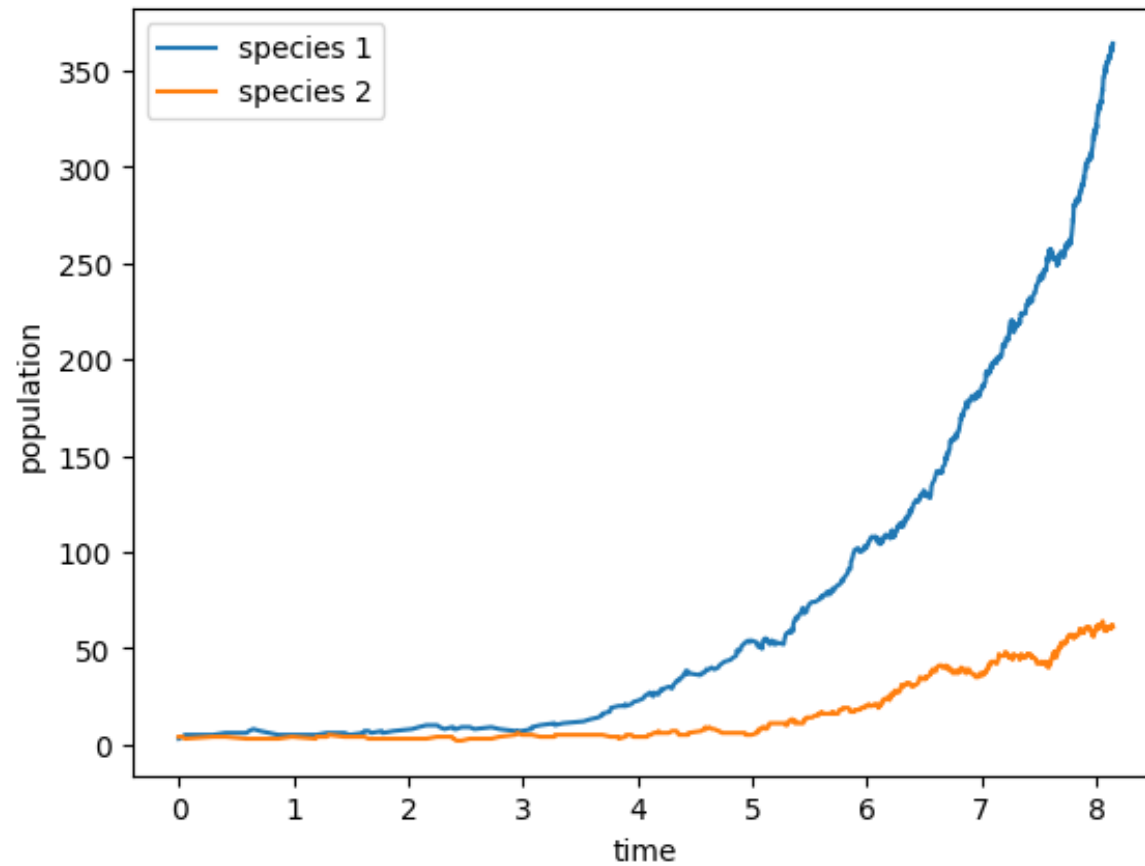
final_T, events = bisse(lam1, lam2, mu1, mu2, p12, p21, ini1, ini2, max_time)
nodes = [(0, ini1, ini2)] + events

# plot each of the species in the same figure
import matplotlib.pyplot as plt
fig, ax = plt.subplots()
ax.plot([node[0] for node in nodes], [node[1] for node in nodes], label='species 1')
ax.plot([node[0] for node in nodes], [node[2] for node in nodes], label='species 2')
ax.set_xlabel('time')
ax.set_ylabel('population')
rates_info = f'lam1={lam1:.2f}, lam2={lam2:.2f}, mu1={mu1:.2f}, mu2={mu2:.2f}, p12={p12:.2f}, p21={p21:.2f}'

```

```
ax.text(0.05, 1.1, rates_info, transform=ax.transAxes, fontsize=9, verticalalignment='top', bbox=dict(boxstyle="
ax.legend()
plt.show()
```

lam1=0.93, lam2=0.66, mu1=0.28, mu2=0.61, p12=0.17, p21=0.34



Function Description

The `generate_data` function simulates the dynamics of two interacting species over multiple iterations and collects the resulting data. This function implements the Birth-Death model with transition rates between the species to generate the data series.

data. This function leverages the Birth-Death model with transition rates between the species to generate the data points.

Parameters

- **num_data_points:** The number of data points to generate.
- **max_lamb_rate:** The maximum value for the birth rates λ_1 and λ_2 .
- **max_mniu_rate:** The maximum value for the death rates μ_1 and μ_2 .
- **max_num_initial_population:** The maximum initial population size for both species.
- **max_time:** The maximum simulation time for each iteration.

Function Dynamics

1. **Initialization:** Two empty lists, `X` and `Y`, are created to store the input parameters and the resulting populations, respectively.
2. **Loop through Data Points:** For each data point:
 - Randomly sample birth rates λ_1 and λ_2 from a uniform distribution between 0 and `max_lamb_rate`.
 - Randomly sample death rates μ_1 and μ_2 from a uniform distribution between 0 and `max_mniu_rate`.
 - Randomly sample transition rates p_{12} and p_{21} from a uniform distribution between 0 and 1.
 - Randomly sample a simulation time `time` from a uniform distribution between 0 and `max_time`.
 - Randomly sample initial populations `ini1` and `ini2` from an integer uniform distribution between 0 and `max_num_initial_population`.
3. **Simulation:** For each set of sampled parameters, the `bisse` function is called to simulate the population dynamics over the sampled time period.
4. **Event Recording:** The populations of species 1 and species 2 at the end of the simulation are recorded. If no events occurred during the simulation, the initial populations are used.
5. **Data Collection:** The sampled parameters and the resulting populations are appended to the lists `X` and `Y`.
6. **Return Values:** The function returns two NumPy arrays `X` and `Y` where `X` contains the input parameters for each data point

c. **Return values:** The function returns two NumPy arrays, `X` and `Y`, where `X` contains the input parameters for each data point and `Y` contains the resulting populations of species 1 and species 2.

```
In [ ]: from tqdm.notebook import tqdm

def generate_data(num_data_points, max_lamb_rate, max_mniu_rate, max_num_initial_population, max_time):
    X = []
    Y = []

    for _ in tqdm(range(num_data_points)):
        lam1, lam2 = np.random.uniform(0, max_lamb_rate, size=2)
        mu1, mu2 = np.random.uniform(0, max_mniu_rate, size=2)
        p12, p21 = np.random.uniform(0, 1, size=2)
        time = np.random.uniform(0, max_time)
        ini1, ini2 = np.random.randint(0, max_num_initial_population, size=2)

        final_T, events = bisse(lam1, lam2, mu1, mu2, p12, p21, ini1, ini2, time)

        _, num_specie1, num_specie2 = (0, ini1, ini2) if len(events) == 0 else events[-1]

        X.append([lam1, lam2, mu1, mu2, p12, p21, ini1, ini2, final_T])
        Y.append([num_specie1, num_specie2])

    return np.array(X), np.array(Y)
```

Data Generation and Splitting

Data Generation

The data generation process involves simulating the dynamics of two interacting species over a large number of iterations using the `generate_data` function. The parameters for this process are as follows:

- **num_data_points:** (64×1250)

- **max_lamb_rate:** 1
- **max_mniu_rate:** 1
- **max_num_initial_population:** 5
- **max_time:** 10

The function `generate_data` is called with these parameters to produce the input data `X` and the corresponding output data `Y`.

Note:

- We consider a small time frame of max time 10 because if sampled lambdas are substantially greater than the sampled mnius, then the growth of the population is exponential and this will become computationally unfeasible.

```
In [ ]: import pandas as pd

num_data_points = 64*1250
max_lamb_rate = 1
max_mniu_rate = 1
max_num_initial_population = 5
max_time = 10

X, Y = generate_data(num_data_points, max_lamb_rate, max_mniu_rate, max_num_initial_population, max_time)

df_X = pd.DataFrame(X, columns=['lam1', 'lam2', 'mu1', 'mu2', 'p12', 'p21', 'ini1', 'ini2', 'final_T'])
df_Y = pd.DataFrame(Y, columns=['num_specie1', 'num_specie2'])
df = pd.concat([df_X, df_Y], axis=1)

from sklearn.model_selection import train_test_split

train_set, test_set = train_test_split(df, test_size=0.1)
train_set, val_set = train_test_split(train_set, test_size=0.2)

0%|          | 0/80000 [00:00<?, ?it/s]
```

```
In [ ]: train_set.corr()
```


Out []:

	lam1	lam2	mu1	mu2	p12	p21	ini1	ini2	final_T	num_speci
lam1	1.000000	-0.000070	0.000534	0.001622	0.003461	0.005749	0.004263	-0.004258	-0.005277	0.205866
lam2	-0.000070	1.000000	-0.010784	-0.000992	0.000831	-0.004449	-0.004942	-0.015801	-0.010984	0.074091
mu1	0.000534	-0.010784	1.000000	0.009937	0.007361	0.000910	0.000710	-0.005827	0.006041	-0.220685
mu2	0.001622	-0.000992	0.009937	1.000000	0.008823	-0.004641	0.005835	-0.007039	0.002614	-0.086814
p12	0.003461	0.000831	0.007361	0.008823	1.000000	0.000222	0.007700	-0.002718	-0.005402	-0.124361
p21	0.005749	-0.004449	0.000910	-0.004641	0.000222	1.000000	0.005385	-0.001673	-0.000859	0.056205
ini1	0.004263	-0.004942	0.000710	0.005835	0.007700	0.005385	1.000000	0.002921	-0.010974	0.077465
ini2	-0.004258	-0.015801	-0.005827	-0.007039	-0.002718	-0.001673	0.002921	1.000000	-0.005296	0.059690
final_T	-0.005277	-0.010984	0.006041	0.002614	-0.005402	-0.000859	-0.010974	-0.005296	1.000000	0.157798
num_specie1	0.205866	0.074091	-0.220685	-0.086814	-0.124361	0.056205	0.077465	0.059690	0.157798	1.000000
num_specie2	0.069280	0.208831	-0.087431	-0.222439	0.047301	-0.123966	0.051084	0.086869	0.153907	0.318691

```

In [ ]: train_X = train_set.drop(columns=['num_specie1', 'num_specie2'])
train_Y = train_set[['num_specie1', 'num_specie2']].copy()

val_X = val_set.drop(columns=['num_specie1', 'num_specie2'])
val_Y = val_set[['num_specie1', 'num_specie2']].copy()

test_X = test_set.drop(columns=['num_specie1', 'num_specie2'])
test_Y = test_set[['num_specie1', 'num_specie2']].copy()

#see missing values
print(train_X.isnull().sum())
print(train_Y.isnull().sum())

print(val_X.isnull().sum())
print(val_Y.isnull().sum())

```

```
print(test_X.isnull().sum())
print(test_Y.isnull().sum())
```

```
lam1      0
lam2      0
mu1       0
mu2       0
p12       0
p21       0
ini1      0
ini2      0
final_T   0
dtype: int64
num_specie1    0
num_specie2    0
dtype: int64
lam1      0
lam2      0
mu1       0
mu2       0
p12       0
p21       0
ini1      0
ini2      0
final_T   0
dtype: int64
num_specie1    0
num_specie2    0
dtype: int64
lam1      0
lam2      0
mu1       0
mu2       0
p12       0
p21       0
ini1      0
```

```

ini1      0
ini2      0
final_T    0
dtype: int64
num_specie1    0
num_specie2    0
dtype: int64

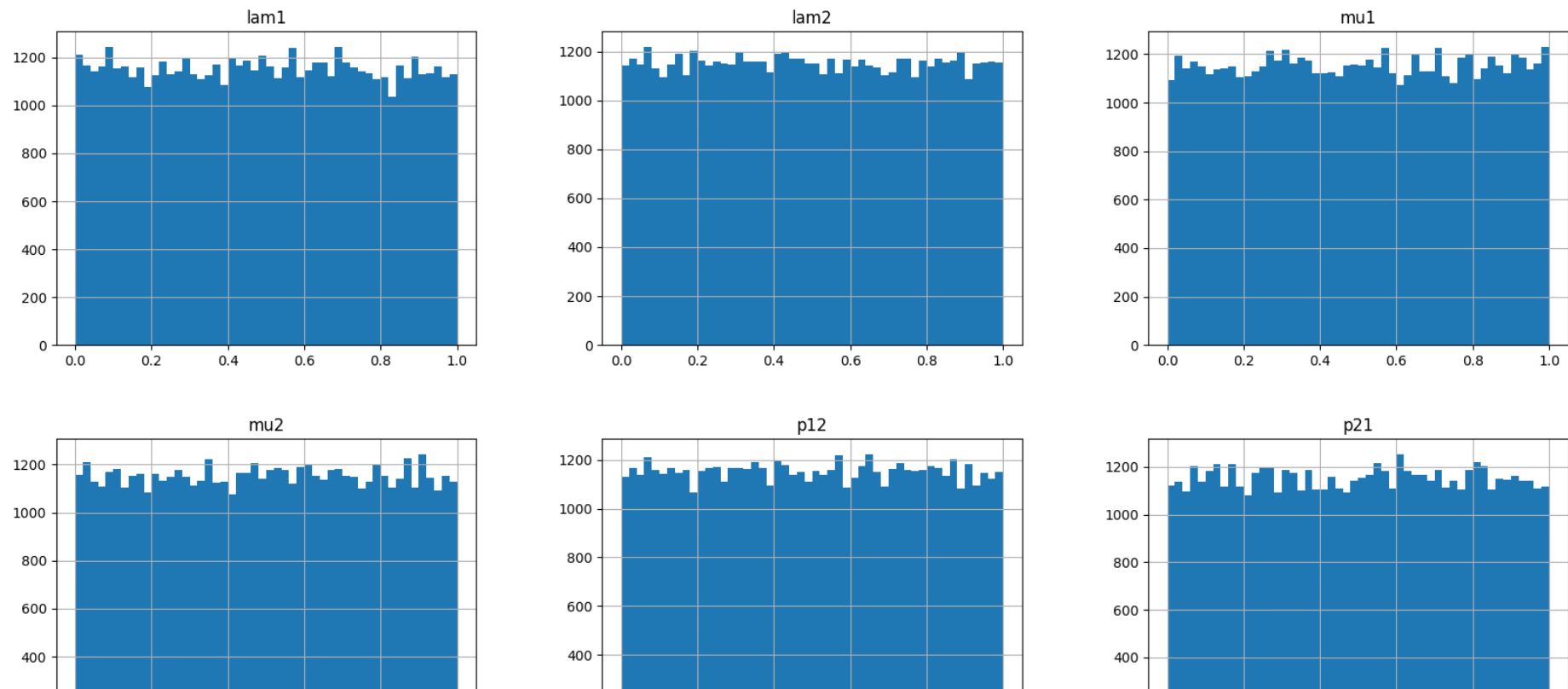
```

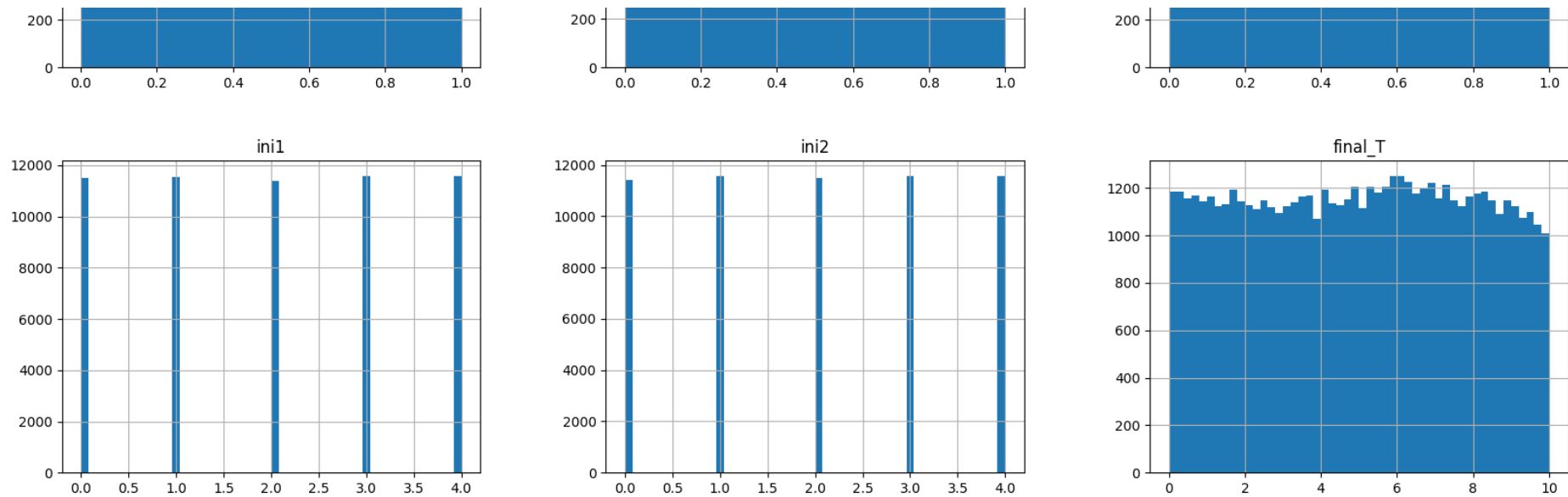
Before we apply any transformation, we should see if distribution of each attribute is heavy tail.

```

In [ ]: import matplotlib.pyplot as plt
train_X.hist(bins=50, figsize=(20, 15))
plt.show()

```





From above we can see the distribution is uniform

Now we can start building the pipeline

```
In [ ]: from sklearn.pipeline import Pipeline
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import MinMaxScaler

pipeline = Pipeline([
    ('imputer', SimpleImputer(strategy='mean')),
    ('min_max_scaler', MinMaxScaler(feature_range=(-1, 1)))
])

train_X_prepared = pipeline.fit_transform(train_X)
```

```
In [ ]: train_X_prepared
```

```
Out[ ]: array([[ 0.85966154,  0.86625492, -0.88649332, ..., -1.          ,
                -0.5          ,  0.79953382],
               [ 0.89198854,  0.18058109,  0.04330034, ...,  0.5          ,
                -1.          , -0.6881473 ],
               [-0.2677207 ,  0.66436848, -0.34219325, ...,  1.          ,
                -0.5          , -0.96937705],
               ...,
               [-0.94178962, -0.08779862,  0.18552385, ..., -0.5          ,
                0.          , -0.09070316],
               [-0.67850873, -0.86402238, -0.60835774, ...,  0.5          ,
                0.5          , -0.73872887],
               [ 0.10871981, -0.81089319, -0.97158014, ..., -0.5          ,
                0.5          , -0.60328665]])
```

```
In [ ]: val_X_prepared = pipeline.transform(val_X)
        test_X_prepared = pipeline.transform(test_X)
        train_Y_prepared = train_Y.to_numpy()
        val_Y_prepared = val_Y.to_numpy()
        test_Y_prepared = test_Y.to_numpy()
```

Neural Network Model Training and Visualization

Neural Network Model Architecture

The code defines a neural network model using TensorFlow's Keras API. The model architecture consists of:

- **Input Layer:** Defined by the shape of `X_train[0]`, which corresponds to the shape of the input data.
- **Dense Layers:** Four hidden layers with 16, 26, 18, and 8 neurons respectively, each using ReLU (Rectified Linear Unit) activation function.
- **Output Layer:** An output layer with neurons equal to the number of outputs (`Y_train[0].shape[0]`), which predicts the populations of species 1 and species 2.

Callback

- **Early Stopping:** A callback (`ea_callback`) is used to monitor the validation loss (`val_loss`). Training will stop early if the validation loss does not improve for 5 consecutive epochs (`patience=5`). The model will restore the weights that give the best validation loss (`restore_best_weights=True`).

Model Compilation

The model is compiled using the Adam optimizer (`optimizer='adam'`) and mean squared error (`loss='mse'`) as the loss function. The accuracy metric is used for evaluation (`metrics=['accuracy']`).

Model Training

The `model.fit` method is called to train the model:

- **X_train, Y_train:** Training data and labels.
- **epochs:** Number of epochs set to 25.
- **batch_size:** Batch size set to 32.
- **validation_data:** Validation data and labels provided as (`X_val, Y_val`).
- **callbacks:** Early stopping callback (`ea_callback`) is passed to monitor validation loss during training.

Training History Visualization

After training, the accuracy and validation accuracy over epochs are plotted using Matplotlib to visualize the model's performance.

```
In [ ]: np.array(train_X_prepared)
```

```
Out [ ]: array([[ 0.85966154,  0.86625492, -0.88649332, ..., -1.          ,
                 -0.5          ,  0.79953382],
                [ 0.89198854,  0.18058109,  0.04330034, ...,  0.5          ,
                 -1.          , -0.6881473 ],
                [-0.2677207 ,  0.66436848, -0.34219325, ...,  1.          ,
                 0.88877357,  0.88877357])
```

```

-0.5          , -0.96937705],
...,
[-0.94178962, -0.08779862,  0.18552385, ..., -0.5          ,
 0.          , -0.09070316],
[-0.67850873, -0.86402238, -0.60835774, ...,  0.5          ,
 0.5          , -0.73872887],
[ 0.10871981, -0.81089319, -0.97158014, ..., -0.5          ,
 0.5          , -0.60328665]])

```

```
In [ ]: train_Y_prepared[0].shape[0]
```

```
Out[ ]: 2
```

```
In [ ]: import tensorflow as tf
```

```

model = tf.keras.Sequential([
    tf.keras.layers.InputLayer(train_X_prepared[0].shape),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(26, activation='relu'),
    tf.keras.layers.Dense(18, activation='relu'),
    tf.keras.layers.Dense(8, activation='relu'),
    tf.keras.layers.Dense(train_Y_prepared[0].shape[0])
])

ea_callback = tf.keras.callbacks.EarlyStopping(monitor='val_loss', patience=5, restore_best_weights=True)

model.compile(optimizer='adam', loss='mse', metrics=['accuracy'])
history = model.fit(train_X_prepared, train_Y, epochs=25, batch_size=32, validation_data=(val_X_prepared, val_Y))

fig, ax = plt.subplots(1, 2, figsize=(15, 5))

ax[0].plot(history.history['loss'], label='train loss')
ax[0].plot(history.history['val_loss'], label='val loss')
ax[0].set_xlabel('epochs')
ax[0].set_ylabel('loss')
ax[0].legend()

```


```

ax[1].plot(history.history['accuracy'], label='train accuracy')
ax[1].plot(history.history['val_accuracy'], label='val accuracy')
ax[1].set_xlabel('epochs')
ax[1].set_ylabel('accuracy')
ax[1].legend()


plt.show()

```


Epoch 1/25

1800/1800  **1s** 404us/step - accuracy: 0.6331 - loss: 2018.6576 - val_accuracy: 0.4796 - val_loss: 1027.3483


Epoch 2/25

1800/1800  **1s** 360us/step - accuracy: 0.4450 - loss: 1072.3293 - val_accuracy: 0.4078 - val_loss: 826.0690


Epoch 3/25

1800/1800  **1s** 369us/step - accuracy: 0.3934 - loss: 918.1927 - val_accuracy: 0.3794 - val_loss: 766.6844


Epoch 4/25

1800/1800  **1s** 357us/step - accuracy: 0.4228 - loss: 825.2649 - val_accuracy: 0.6181 - val_loss: 736.8059


Epoch 5/25

1800/1800  **1s** 355us/step - accuracy: 0.5710 - loss: 749.8409 - val_accuracy: 0.6532 - val_loss: 714.8970


Epoch 6/25

1800/1800  **1s** 357us/step - accuracy: 0.6494 - loss: 768.8055 - val_accuracy: 0.7015 - val_loss: 714.7110


Epoch 7/25

1800/1800  **1s** 358us/step - accuracy: 0.6748 - loss: 726.2449 - val_accuracy: 0.7148 - val_loss: 712.2497

Epoch 8/25

1800/1800  **1s** 372us/step - accuracy: 0.6778 - loss: 713.0627 - val_accuracy: 0.6917 - val_loss: 651.6809

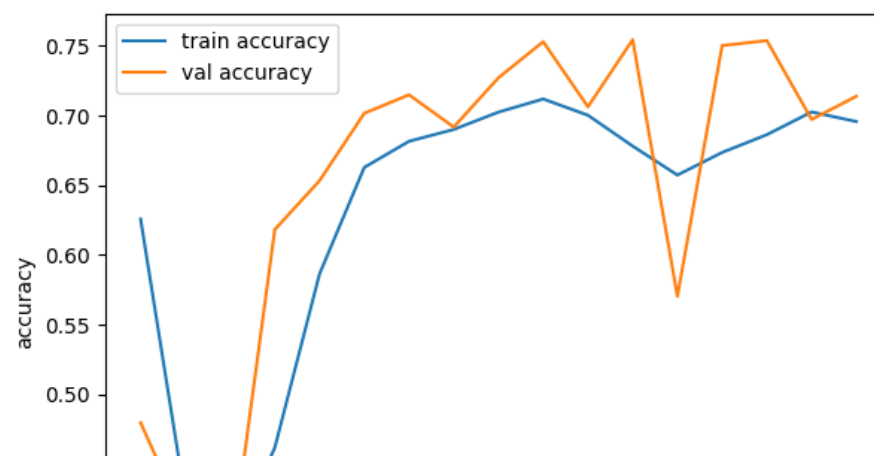
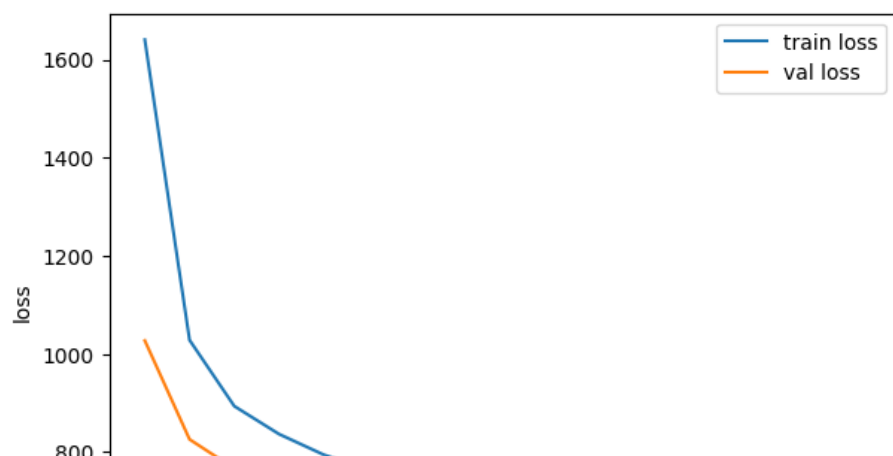
Epoch 9/25

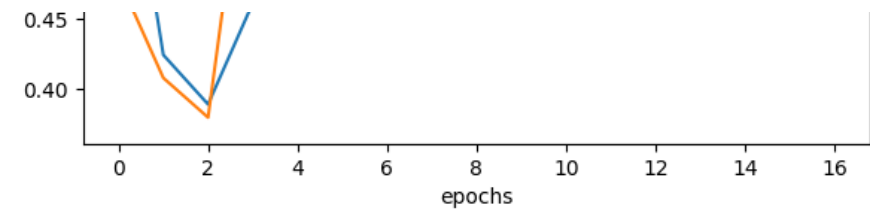
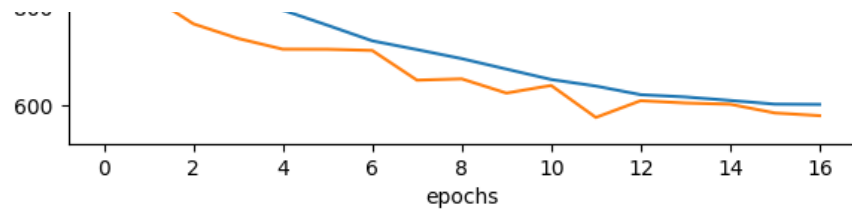
1800/1800  **1s** 365us/step - accuracy: 0.7146 - loss: 690.7686 - val_accuracy: 0.7270 - val_loss: 654.3060

Epoch 10/25

1800/1800 ————— **1s** 361us/step – accuracy: 0.6996 – loss: 682.2421 – val_accuracy: 0.7529 – val_loss: 625.2794
 Epoch 11/25
1800/1800 ————— **1s** 361us/step – accuracy: 0.7089 – loss: 662.0202 – val_accuracy: 0.7062 – val_loss: 640.8795
 Epoch 12/25
1800/1800 ————— **1s** 360us/step – accuracy: 0.6768 – loss: 654.2625 – val_accuracy: 0.7544 – val_loss: 575.5341
 Epoch 13/25
1800/1800 ————— **1s** 374us/step – accuracy: 0.6676 – loss: 649.3317 – val_accuracy: 0.5701 – val_loss: 609.7438
 Epoch 14/25
1800/1800 ————— **1s** 358us/step – accuracy: 0.6701 – loss: 631.6699 – val_accuracy: 0.7501 – val_loss: 605.0446
 Epoch 15/25
1800/1800 ————— **1s** 358us/step – accuracy: 0.6861 – loss: 598.2549 – val_accuracy: 0.7537 – val_loss: 602.5328
 Epoch 16/25

1800/1800 ————— **1s** 356us/step – accuracy: 0.7051 – loss: 587.5711 – val_accuracy: 0.6971 – val_loss: 584.7032
 Epoch 17/25
1800/1800 ————— **1s** 369us/step – accuracy: 0.6952 – loss: 597.6011 – val_accuracy: 0.7138 – val_loss: 579.2067





Model Evaluation on Test Data

Model Evaluation

To evaluate the trained neural network model on the test data (`X_test` and `Y_test`), the `model.evaluate` method is used. This method computes the loss and metrics (accuracy in this case) on the test set.

```
In [ ]: cost, acc = model.evaluate(test_X_prepared, test_Y_prepared)
print(f'Test accuracy: {acc:.3f}')
```

250/250 ————— 0s 259us/step – accuracy: 0.7518 – loss: 615.5843
Test accuracy: 0.753

Saving the model

We can observe that our current model has an accuracy of 72.2% on the test set. In other words, given a vector space of parameters within the bounds previously defined, we can predict the number of of each species with an accuracy of 72.2%. We can proceed and save the model for future use.

```
In [ ]: # save the model
model.save('bisse_model.keras')
```