

Academic Year: 2023-24

Semester: VI

Class / Branch: TE-IT

Subject: BI Lab

Name of Instructor: Prof. Apeksha Mohite

Name of Student: Aashay Ingale

Student ID: 21104006

Date of Performance: 30/03/2023 Date of Submission: 30/03/2023

Experiment No. 13

Aim: Business Intelligence Mini Project.

1. Problem Definition: Will Construction of Motorway harm the nearby Amphibians.

2. Data mining task to be performed:

The data mining task performed is clustering. Specifically, the two algorithms mentioned are Simple K Means and EM (Expectation-Maximization). Both algorithms are commonly used for clustering tasks in data mining.

- a. **Simple K Means (K-Means):** This algorithm aims to partition the dataset into K clusters, where each instance belongs to the cluster with the nearest mean. It iteratively assigns instances to clusters based on the Euclidean distance between the instance and cluster centroids, then updates the centroids based on the mean of instances in each cluster. In the provided information, Simple K Means was applied to cluster the dataset into two clusters.
- b. **EM** (**Expectation-Maximization**): EM algorithm is a probabilistic approach to clustering that assumes the data is generated from a mixture of several Gaussian distributions. It iteratively estimates the parameters of these Gaussian distributions (mean and variance) to maximize the likelihood of the observed data. EM algorithm is particularly useful when dealing with data that may have overlapping clusters or clusters with different shapes and sizes.

Therefore, the data mining tasks performed in the provided information are indeed Simple K Means and EM clustering algorithms.



3. Dataset identified: Amphibians

4. Source of dataset: https://archive.ics.uci.edu/dataset/528/amphibians

5. Details of the dataset:

The dataset is a multilabel classification problem. The goal is to predict the presence of amphibian's species near the water reservoirs based on features obtained from GIS systems and satellite images.

The information we're talking about comes from maps, satellite images, and studies done to see how building new roads might affect nature, specifically amphibians like frogs and salamanders, in two places in Poland.

For Road A, which is part of a big motorway plan, researchers checked out a stretch of land about 500 meters wide on each side of where the road would go. They did this in 2010 and 2011, and also later from 2014 to 2016. They found about 80 places where amphibians were likely to lay eggs or live.

For Road B, which is part of another motorway plan, they looked at two different routes in a certain area. They looked at maps, old data, and went out into the field to see where amphibians were. They did this all-in springtime. They found about 125 places where amphibians lived or might live. They did a similar thing with the land about 500 meters wide on each side of the road paths. This time, they found 109 spots where amphibians were likely to be.



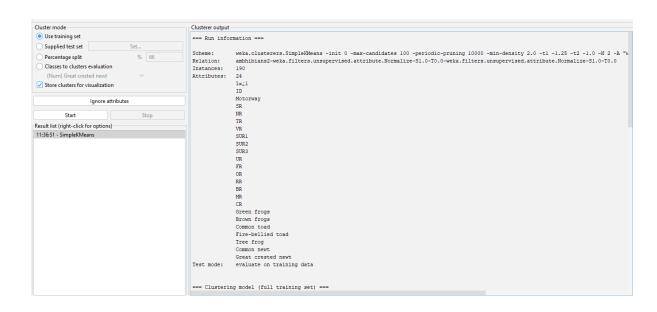


6. Algorithms to accomplish the task:

Clustering:

a. SimplekMeans

The Simple K Means clustering algorithm was applied to a dataset consisting of 190 instances with 24 attributes. This clustering process yielded two distinct clusters, where Cluster 0 comprises 110 instances, representing approximately 58% of the dataset, while Cluster 1 contains 80 instances, accounting for around 42% of the data. Each cluster demonstrates unique characteristics as illustrated by their centroids, which act as representative points for the cluster. Cluster 0's centroid exhibits relatively lower values for attributes such as 'SR', 'NR', 'TR', and 'Green frogs' compared to Cluster 1, suggesting a different profile for this group. Conversely, Cluster 1 displays higher values for attributes like 'SR', 'NR', 'FR', 'OR', 'Green frogs', 'Brown frogs', 'Common toad', and 'Common newt', indicating distinct patterns within this cluster. This analysis underscores the heterogeneous nature of the dataset and furnishes valuable insights for further exploration and decisionmaking processes.





Parshvanath Charitable Trust's .. P. SHAH INSTITUTE OF TECHNOLOGY, THANE (All Programs Accredited by NBA)



Department of Information Technology ----- Number of iterations: 5 Within cluster sum of squared errors: 461.2770383033899

Number of iterations Within cluster sum of		rors: 461.27	70383033899		
Initial starting pos	ints (random):	:			
				5,0.714286,0.909091,0.181818,0,0,0.5,0.2,0.5,0,0.5,0,1,0,0,0,0,0 67,0.25,0.071429,0.909091,0.181818,1,0.75,1,0,0.1,0,0.5,1,1,1,0,1,0,0	
Missing values globa	ally replaced	with mean/m	ode		
Final cluster centro	oids:				
		Cluster#			
Attribute	Full Data	0	1 (00.0)		
	(190.0)		(80.0)		
l		0.3708			
ID		0.3675			
Motorway		A1			
SR	0.0192	0.006	0.0373		
NR	0.1298	0.1091	0.1583		
TR	0.3284	0.4709			
VR	0.4737	0.5159	0.4156		
SUR1	0.3008	0.2773	0.333		
SUR2		0.4372			
SUR3		0.5248			
UR		0.1848			
FR		0.1386			
OR		0.9021			
RR			0.1875		
BR		0.3045			
MR	0.0237		0.0063		
CR	0.5053	0.5045	0.5062		
_					
CR	0.5053		0.5062		
Green frogs			0.8875		
Brown frogs Common toad		0.6545 0.5182			
Fire-bellied toad		0.1273	0.55		
Tree frog		0.1273			
Common newt		0.2273			
Great crested newt			0.2		
Time taken to build	model (full t	training dat	a) • 0 02 s	econds	
=== Model and evalua	ation on train	ning set ===			
Clustered Instances					
0 110 (58%)					
1 80 (42%)					
20 (120)					
				Lo	a

b. EM

- i. Clustering Model: The Expectation-Maximization (EM) algorithm was used to cluster the dataset. The algorithm was configured with various parameters such as the number of iterations, maximum clusters, and convergence criteria.
- ii. Number of Clusters: The EM algorithm determined that there are five clusters in the dataset.



iii. Cluster Information:

- Cluster 0: 55 instances (29%)
- Cluster 1: 36 instances (19%)
- Cluster 2: 30 instances (16%)
- Cluster 3: 8 instances (4%)
- Cluster 4: 61 instances (32%)
- **iv. Cluster Centroids:** For each cluster, the output provides the mean and standard deviation of attribute values. Attributes include various features such as motorway type, surface type, and the presence of different amphibian species. These centroid values give insights into the typical characteristics of instances within each cluster.

v. For example:

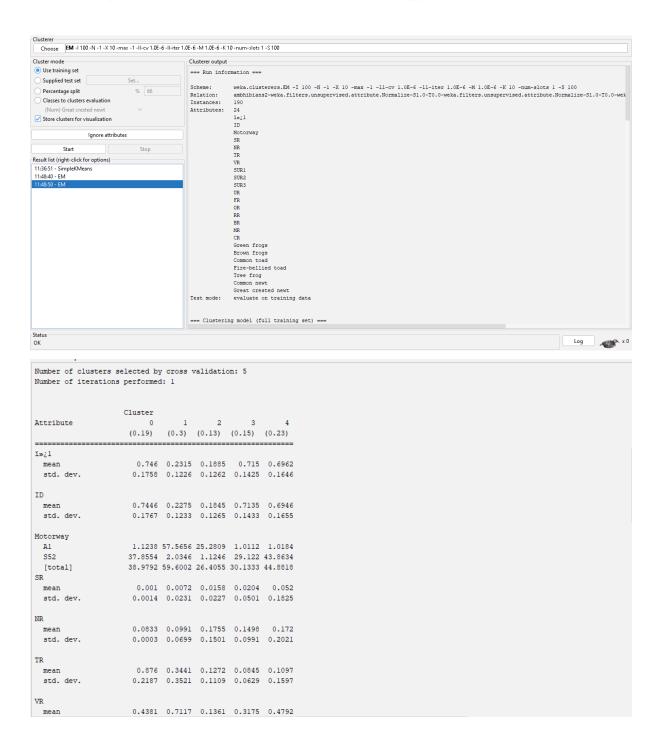
- Cluster 0 has higher mean values for attributes such as "Motorway A1", "Green frogs", "Common toad", and "Common newt".
- Cluster 4 has higher mean values for attributes such as "Motorway "S52", "Brown frogs", "Fire-bellied toad", and "Great crested newt".
- vi. Log Likelihood: The log likelihood value of 17.50416 indicates how well the model fits the data. Higher log likelihood values generally indicate better model fit.
- vii. In summary, the EM algorithm clustered the dataset into five distinct groups based on the attributes provided. Analysis of cluster centroids provides insights into the characteristics of each cluster, while the log likelihood value assesses the overall quality of the clustering model.



Parshvanath Charitable Trust's A. P. SHAH INSTITUTE OF TECHNOLOGY, THANE (All Programs Accredited by NBA)



Department of Information Technology





A. P. SHAH INSTITUTE OF TECHNOLOGY, THANE (All Programs Accredited by NBA)



Department of Information Technology

SUR1					
mean	0.2751				
std. dev.	0.2707	0.2073	0.2196	0.2509	0.2657
SUR2					
	0.5035	0.3804	0.5213	0.5488	0.5583
std. dev.	0.348				
SUR3					
mean	0.5582				
std. dev.	0.3017	0.3237	0.2974	0.2464	0.29
UR					
mean	0.009	0.0156	0.8694	0.5307	0.3642
std. dev.					
FR					
mean			0.6363		
std. dev.	0.0004	0.0565	0.3281	0.3909	0.3334
OR					
mean	0.9524	0.9327	0.7586	0.7819	0.9463
std. dev.					
RR					
mean	0.2459				
std. dev.	0.2548	0.2885	0.2244	0.2272	0.1764
RD					
BR mean	0.2856	0.3787	0.1439	0.1236	0.185
std. dev.					
	0.200				
MR					
mean	0.0271	0.0521	0	0.0178	0
CR					
	0.5	0.5	0.5	0.5355	0.5
	0.5 0.0022	0.5	0.5 0.1431	0.5355 0.1285	0.5
	0.5 0.0022	0.5 0.0012	0.5 0.1431	0.5355 0.1285	0.5
	0.5 0.0022	0.5 0.0012	0.5 0.1431	0.5355 0.1285	0.5
mean std. dev.	0.0022	0.0012	0.1431	0.1285	0.9304
mean std. dev. Green frogs	0.0022	0.0012	0.1431	0.1285	0.9304
mean std. dev. Green frogs mean std. dev.	0.0022	0.0012	0.1431	0.1285	0.9304
mean std. dev. Green frogs mean std. dev. Brown frogs	0.0022 0.1306 0.337	0.0012 0.4189 0.4934	0.1431 0.6145 0.4867	0.1285 0.8582 0.3488	0.9304 0.2545
mean std. dev. Green frogs mean std. dev. Brown frogs mean	0.0022 0.1306 0.337	0.0012 0.4189 0.4934	0.1431 0.6145 0.4867	0.1285 0.8582 0.3488	0.9304 0.2545
mean std. dev. Green frogs mean std. dev. Brown frogs	0.0022	0.0012 0.4189 0.4934	0.1431 0.6145 0.4867	0.1285 0.8582 0.3488	0.9304 0.2545
mean std. dev. Green frogs mean std. dev. Brown frogs mean	0.0022 0.1306 0.337 0.8377 0.3688	0.0012 0.4189 0.4934 0.6751 0.4683	0.1431 0.6145 0.4867 0.4966 0.5	0.1285 0.8582 0.3488 0.8223 0.3823	0.9304 0.2545 1 0.0008
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev.	0.0022 0.1306 0.337	0.0012 0.4189 0.4934 0.6751 0.4683	0.1431 0.6145 0.4867 0.4966 0.5	0.1285 0.8582 0.3488 0.8223 0.3823	0.9304 0.2545 1 0.0008
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad	0.0022 0.1306 0.337 0.8377 0.3688	0.0012 0.4189 0.4934 0.6751 0.4683	0.1431 0.6145 0.4867 0.4966 0.5	0.1285 0.8582 0.3488 0.8223 0.3823	0.9304 0.2545 1 0.0008
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev.	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833	0.0012 0.4189 0.4934 0.6751 0.4683	0.1431 0.6145 0.4867 0.4966 0.5	0.1285 0.8582 0.3488 0.8223 0.3823	0.9304 0.2545 1 0.0008
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833	0.4189 0.4934 0.6751 0.4683 0.5554 0.4969	0.1431 0.6145 0.4867 0.4966 0.5	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478	0.9304 0.2545 1 0.0008 0.9781 0.1465
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833	0.4189 0.4934 0.6751 0.4683 0.5554 0.4969	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478	0.9304 0.2545 1 0.0008 0.9781 0.1465
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833	0.4189 0.4934 0.6751 0.4683 0.5554 0.4969	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478	0.9304 0.2545 1 0.0008 0.9781 0.1465
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev.	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833	0.4189 0.4934 0.6751 0.4683 0.5554 0.4969	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478	0.9304 0.2545 1 0.0008 0.9781 0.1465
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856 0.1302 0.3365	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean std. dev.	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856 0.1302 0.3365 0.0524 0.2228	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean std. dev. Great crested news	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856 0.1302 0.3365	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766 0.174 0.3791	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851 0.1144	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean std. dev. Great crested newt mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856 0.1302 0.3365 0.0524 0.2228	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766 0.174 0.3791	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287 0.1232 0.3285	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851 0.114 0.3178	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295 0.9289 0.2571
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean std. dev. Great crested news	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856 0.1302 0.3365 0.0524 0.2228	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766 0.174 0.3791	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287 0.1232 0.3285	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851 0.114 0.3178	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295 0.9289 0.2571
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean std. dev. Great crested newt mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856 0.1302 0.3365 0.0524 0.2228	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766 0.174 0.3791	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287 0.1232 0.3285	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851 0.114 0.3178	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295 0.9289 0.2571
mean std. dev. Green frogs mean std. dev. Brown frogs mean std. dev. Common toad mean std. dev. Fire-bellied toad mean std. dev. Tree frog mean std. dev. Common newt mean std. dev. Great crested newt mean	0.0022 0.1306 0.337 0.8377 0.3688 0.3719 0.4833 0.1817 0.3856 0.1302 0.3365 0.0524 0.2228	0.0012 0.4189 0.4934 0.6751 0.4683 0.5554 0.4969 0.2087 0.4064 0.3489 0.4766 0.174 0.3791	0.1431 0.6145 0.4867 0.4966 0.5 0.7426 0.4372 0.205 0.4037 0.1232 0.3287 0.1232 0.3285	0.1285 0.8582 0.3488 0.8223 0.3823 0.6466 0.478 0.0214 0.1448 0.379 0.4851 0.114 0.3178	0.9304 0.2545 1 0.0008 0.9781 0.1465 0.7848 0.4109 0.756 0.4295 0.9289 0.2571



Parshvanath Charitable Trust's A. P. SHAH INSTITUTE OF TECHNOLOGY, THANE (All Programs Accredited by NBA)



Department of Information Technology

Fire-bellie	d toad					
mean		0.1817	0.2087	0.205	0.0214	0.7848
std. dev.		0.3856	0.4064	0.4037	0.1448	0.4109
Tree frog						
mean			0.3489			
std. dev.		0.3365	0.4766	0.3287	0.4851	0.4295
Common newt						
std. dev.		0.2228	0.3791	0.3285	0.3178	0.2571
Great crest						
std. dev.		0.0183	0.2812	0.0178	0.0235	0.4835
Time taken	to build mo	del (ful	l traini	ng data)	: 0.61	seconds
=== Model a	nd evaluati	on on tr	aining s	et ===		
Clustered I	nstances					
0 55	(29%)					
1 36	(19%)					
2 30	(16%)					
3 8	(4%)					
4 61	(32%)					
Log likelih	ood: 17.504	16				

7. Conclusion: Thus, we can conclude that these studies helped us understand where amphibians were living so that road construction could be planned in a way that would cause the least harm to them.