

DTSA 5510 - Unsupervised Algorithms in Machine Learning Final Project: Unsupervised Learning in the Case of UAV Image Classification

Israel Johnson
August 21, 2023



University of Colorado **Boulder**

PROBLEM STATEMENT

- There has been ongoing research in the area of applying artificial intelligence and machine learning (AI/ML) techniques in detecting unmanned aerial vehicles (UAVs). These applications usually relate to a combat scenario or military use case.
- Researchers still explore different ML model approaches for analyzing images and video data related to UAVs, especially when images or videos can contain other objects in the environment that can make UAV detection more challenging.

PROJECT OBJECTIVE

- The goal for this project is to apply multiple unsupervised learning model techniques for UAV classification, conduct comparative analysis to find the best unsupervised model within the scope of the project, and make final comparisons of the best unsupervised model results with the results of a supervised learning model approach (in order to determine which specific approach performs better for this use case).
- Three unsupervised learning models are applied to this project, which includes the following: a hierarchical clustering model and two anomaly detection models: Non-negative matrix factorization (NMF) and Multi-layer Perceptron Regressor (MLP). A support vector machine (SVM) classifier will serve as the supervised learning model approach for comparison purposes.

DATA SOURCES

- 1.) Gareth. 2022. *UAV Detection Dataset*. Retrieved from: <https://www.kaggle.com/datasets/nelyg8002000/uav-detection-dataset-images>
- 2.) Ailurophile. 2019. *200 Bird Species with 11,788 Images*. Retrieved from: <https://www.kaggle.com/datasets/veeralakrishna/200-bird-species-with-11788-images>
- 3.) Kaggle (Research Code Competition). 2019. *Generative Dog Images*. Retrieved from: <https://www.kaggle.com/competitions/generative-dog-images/overview>

GITHUB REPOSITORY SOURCE

- **Source:** https://github.com/IsraelsLibrary/DTSA_5510_Unsupervised_Algorithms__In_Machine_Learning



DATA CLEANING PROCESS

- For initial data processing, the software reads in three datasets, each dataset focusing on a different category of images (UAVs, Birds, and Dogs).
- Two helper functions clean the datasets by identifying and removing any corrupted images in each dataset.
- Principal Component Analysis (PCA) and flattening the image data helps in reducing dimensionality of the datasets and processing the data for model analysis.



University of Colorado **Boulder**

DATA CLEANING PROCESS

```
1  # Helper functions that check for any corrupted images in the selected directories.
2  # If present, those corrupted image files are removed.
3
4  def is_corrupted(image_path):
5      try:
6          Image.open(image_path).verify()
7          return False
8      except (IOError, SyntaxError):
9          return True
10
11 def delete_corrupted_images(directory):
12     return [
13         filename
14         for filename in os.listdir(directory)
15         if is_corrupted(os.path.join(directory, filename)) and os.remove(os.path.join(directory, filename))
16     ]
17
```



DATA TRANSFORMATIONS

```
1 # Reading in the image datasets
2 uav_imgs = []
3 birds_imgs = []
4 dogs_imgs = []
5
6 uav_path = "data/uavs_datasets/"
7 bird_path = "data/birds_datasets/"
8 dogs_path = "data/dogs_datasets/"
9 collection = {uav_path:uav_imgs, bird_path:birds_imgs, dogs_path:dogs_imgs}
10
11 # Checking and removing any corrupted images from the file directories. Storing the valid images
12 # into a temporary data structure, with image datasets assigned to their corresponding image categories.
13 for path, imgs in collection.items():
14     # Remove corrupted images before loading
15     corrupted_images = delete_corrupted_images(path)
16     if len(corrupted_images) > 0:
17         print("Corrupted images removed:", corrupted_images)
18     else:
19         print("Found no corrupted images in the selected directory.")
20
21     goodies = [img for img in os.listdir(path) if img.endswith('.jpg')]
22     imgs.extend([np.array(Image.open(os.path.join(path, img))) for img in goodies])
```

```
1 # Reformatting the custom dataset to contain an equal number of images
2 # for each category (for further testing)
3 mod_collection = {'UAVs': collection["data/uavs_datasets/"][:5],
4                  'Birds': collection["data/birds_datasets/"][:5],
5                  'Dogs': collection["data/dogs_datasets/"][:5]}
6
7 df = pd.DataFrame(mod_collection)
```



DATA TRANSFORMATIONS

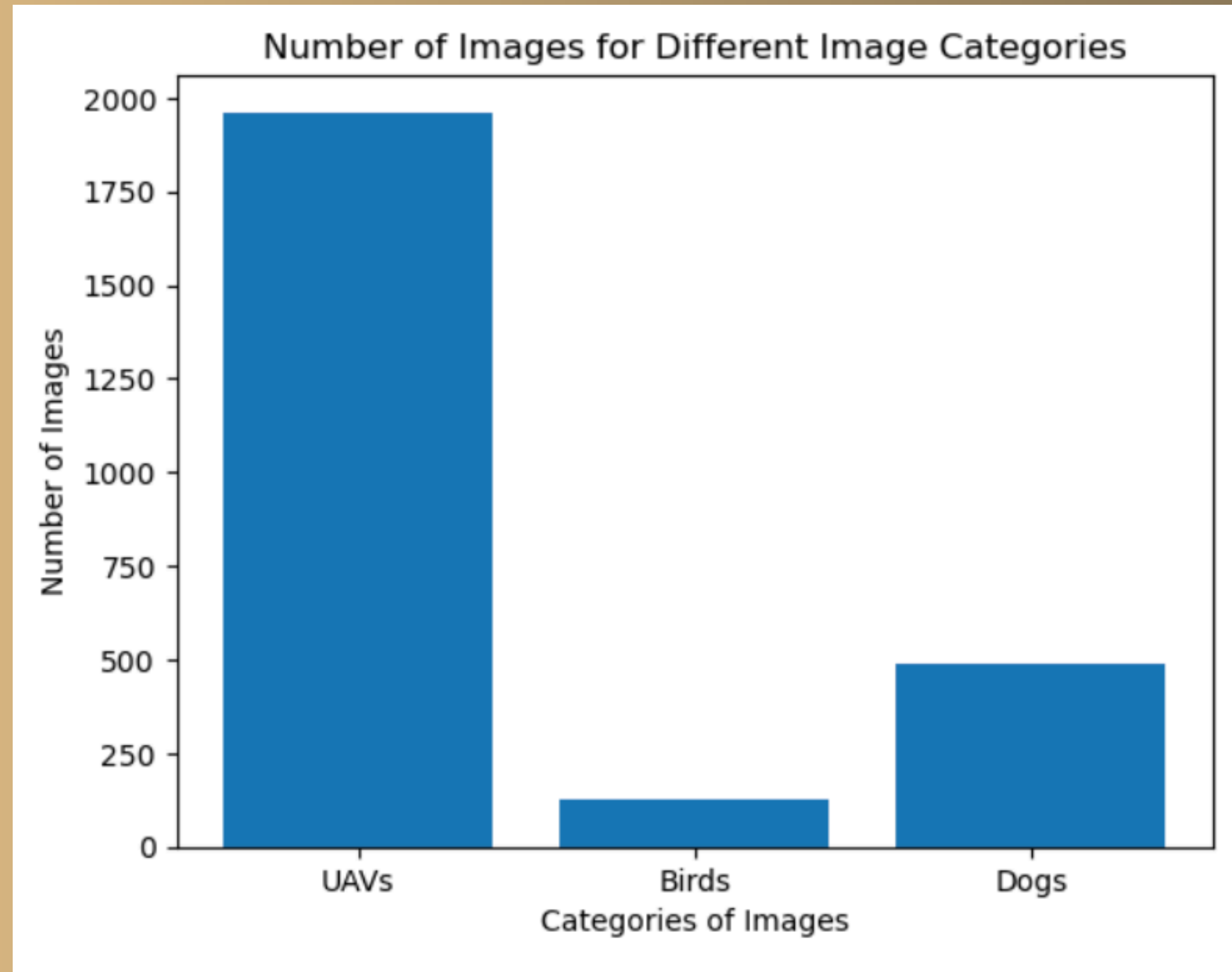
```
1 # Step 1: Data Preprocessing
2
3 # Processing and flattening the input matrices in preparation for model training
4
5 # Dividing the data up by categories from the custom dataframe.
6 uav_pixel_data = df['UAVs'].values
7 bird_pixel_data = df['Birds'].values
8 dog_pixel_data = df['Dogs'].values
9
10 # Flattening the matrix values from the datasets by category
11 uav_flattened_data = [pixel for image in uav_pixel_data for row in image for pixel in row]
12 bird_flattened_data = [pixel for image in bird_pixel_data for row in image for pixel in row]
13 dog_flattened_data = [pixel for image in dog_pixel_data for row in image for pixel in row]
14
15 # Converting the flattened data into feature matrices
16 uav_feature_matrix = np.array(uav_flattened_data)
17 bird_feature_matrix = np.array(bird_flattened_data)
18 dog_feature_matrix = np.array(dog_flattened_data)
19
20 # Applying PCA for dimensionality reduction
21 num_components = 3
22 pca = PCA(n_components=num_components)
23
24 uav_preprocessed_data = pca.fit_transform(uav_feature_matrix[:100])
25 bird_preprocessed_data = pca.fit_transform(bird_feature_matrix[:100])
26 dog_preprocessed_data = pca.fit_transform(dog_feature_matrix[:100])
27
28 preprocessed_data = np.vstack((uav_preprocessed_data, bird_preprocessed_data, dog_preprocessed_data))
29
```



EXPLORATORY DATA ANALYSIS (EDA)

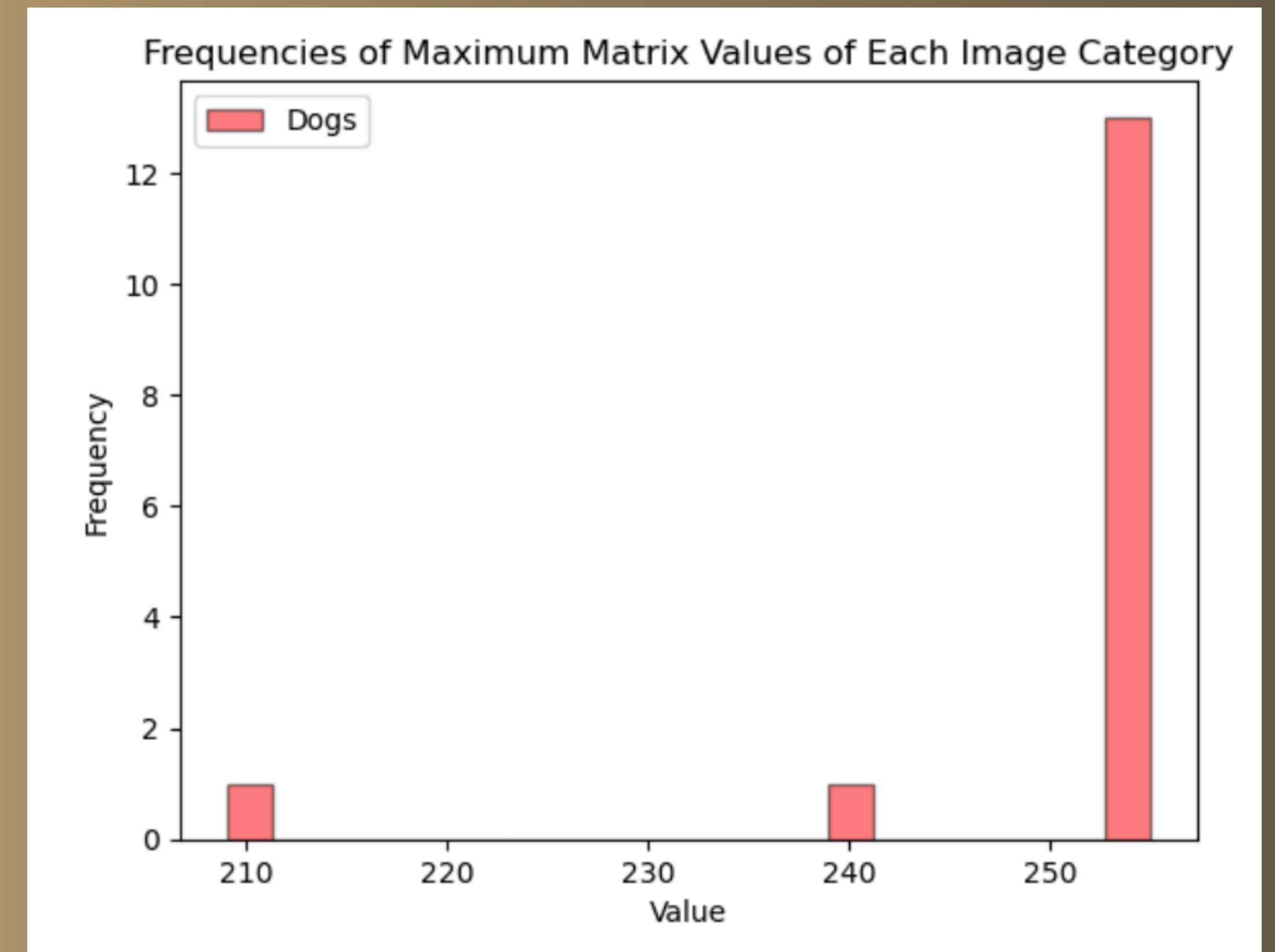
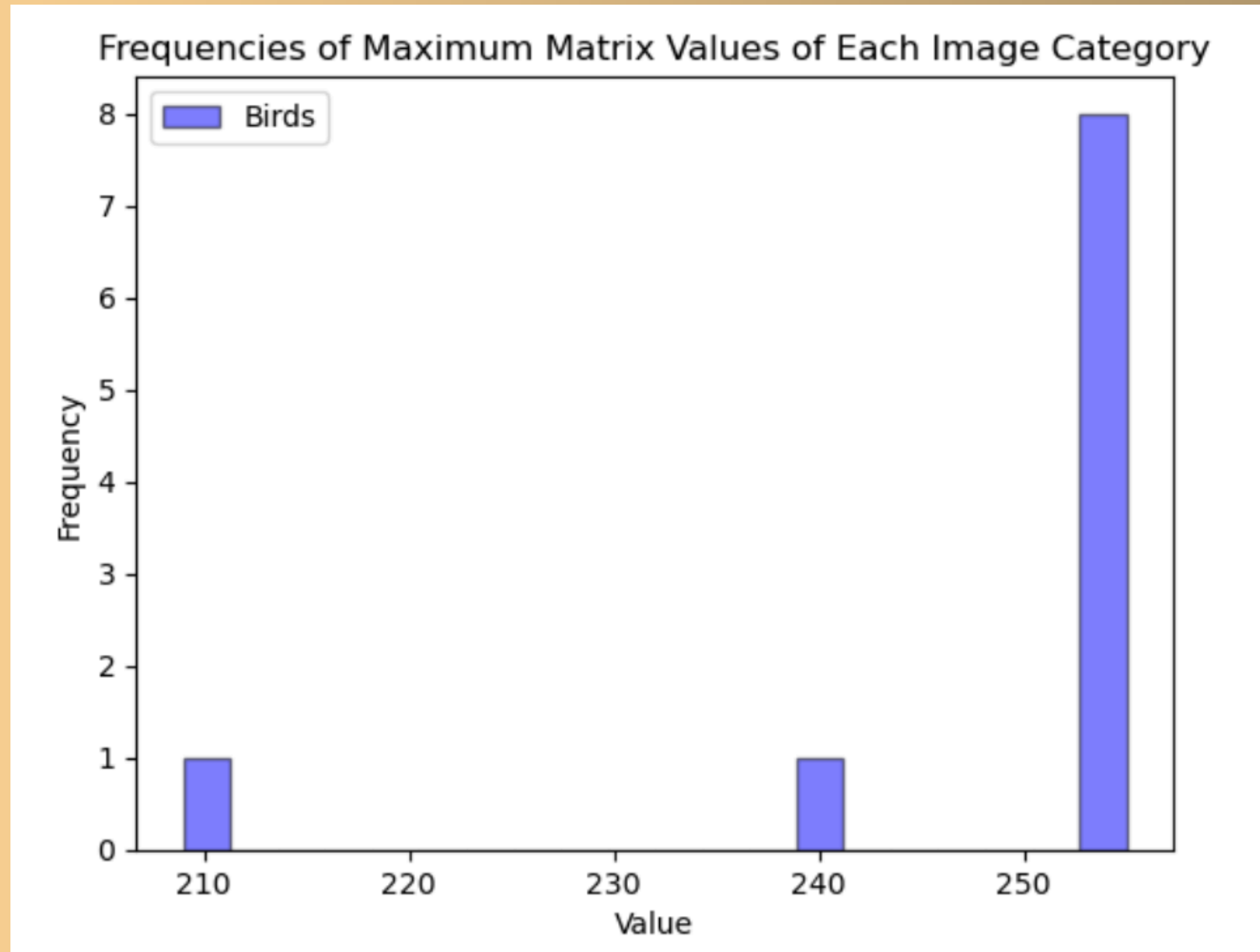
- Generated histograms to track the frequencies of the minimum values, maximum values, and average values of image matrices for each category of images.
- The histograms revealed enough differences between image data frequencies to confirm that there is a sufficient amount of project data to perform image classification.

EXPLORATORY DATA ANALYSIS

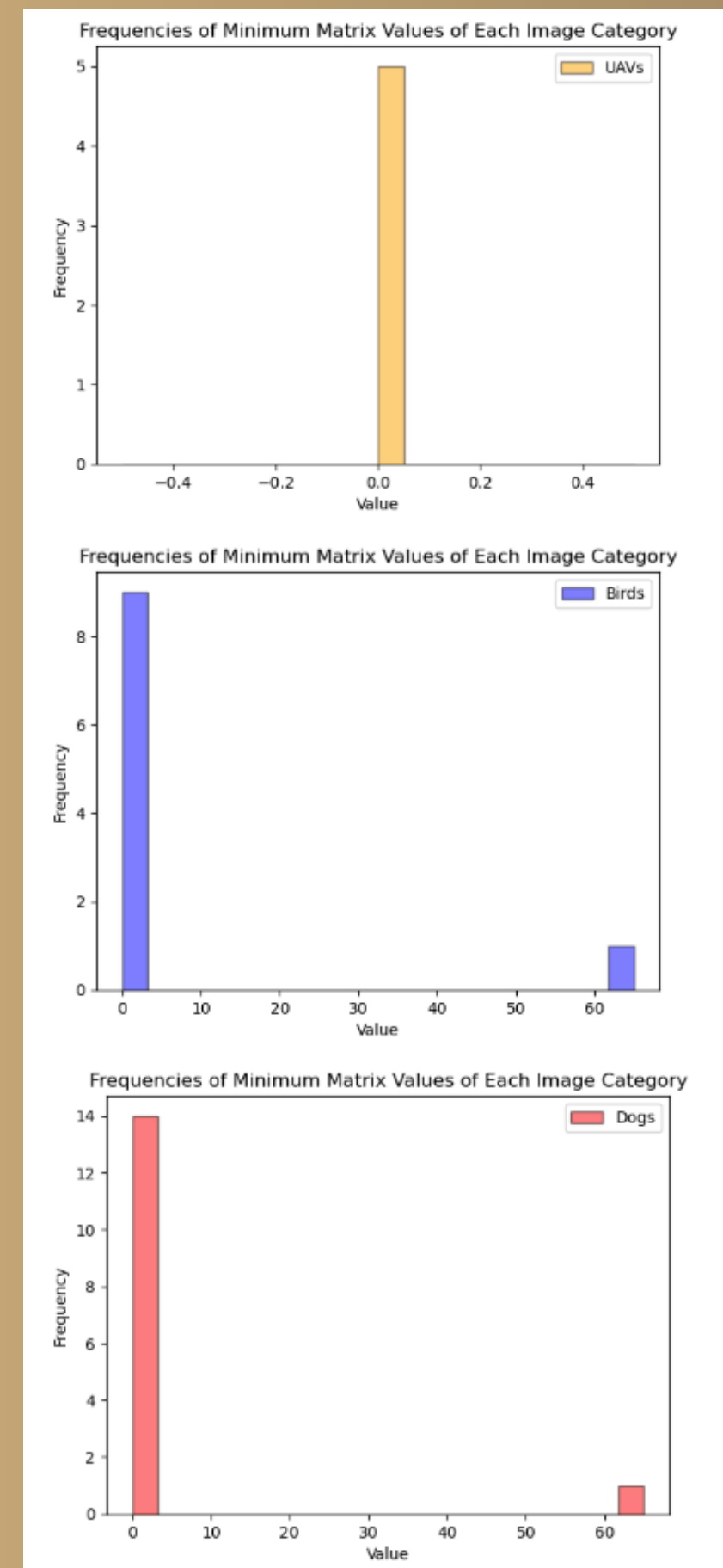


University of Colorado **Boulder**

EXPLORATORY DATA ANALYSIS

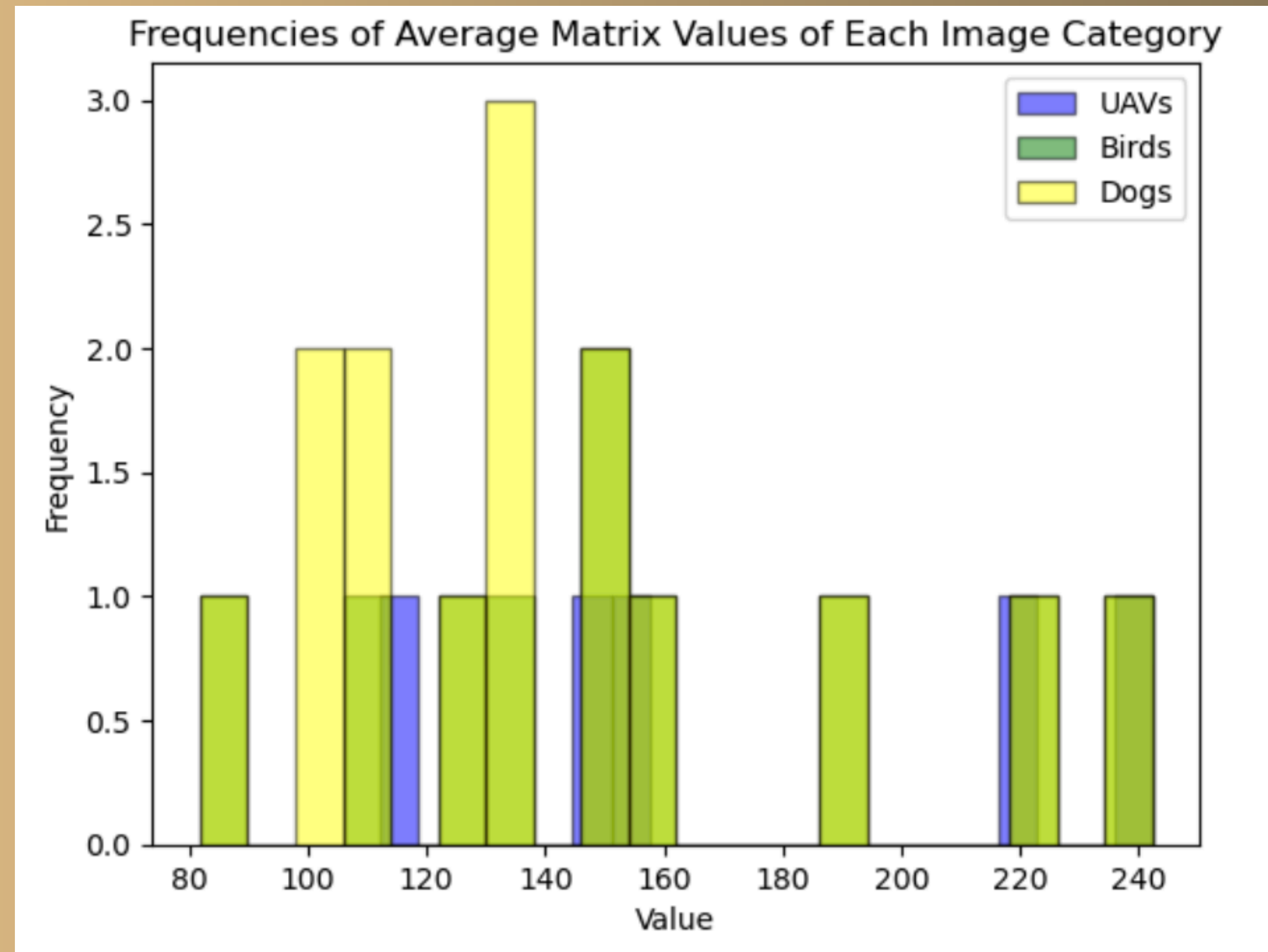


EXPLORATORY DATA ANALYSIS



University of Colorado **Boulder**

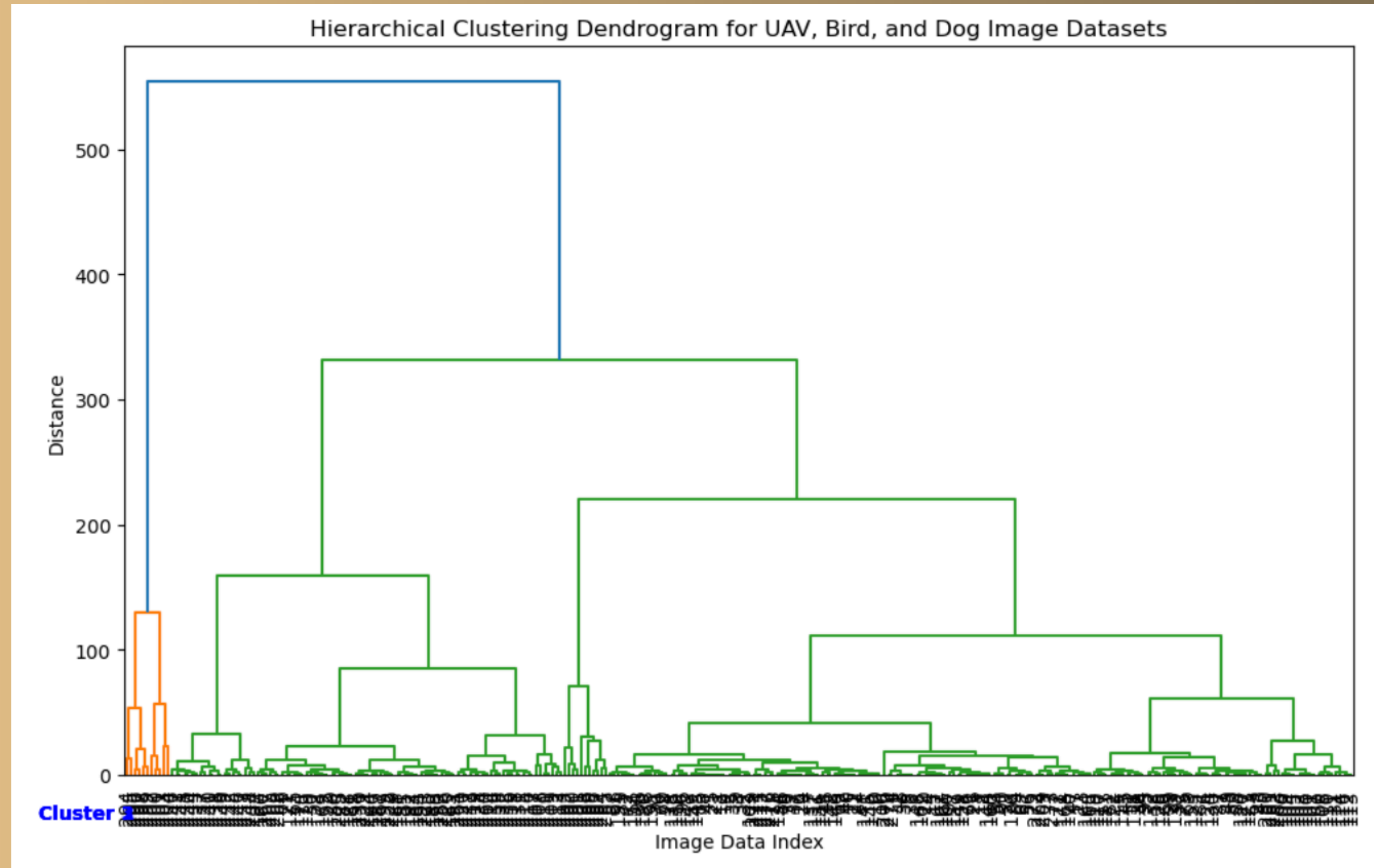
EXPLORATORY DATA ANALYSIS



MODEL APPROACH AND ANALYSIS

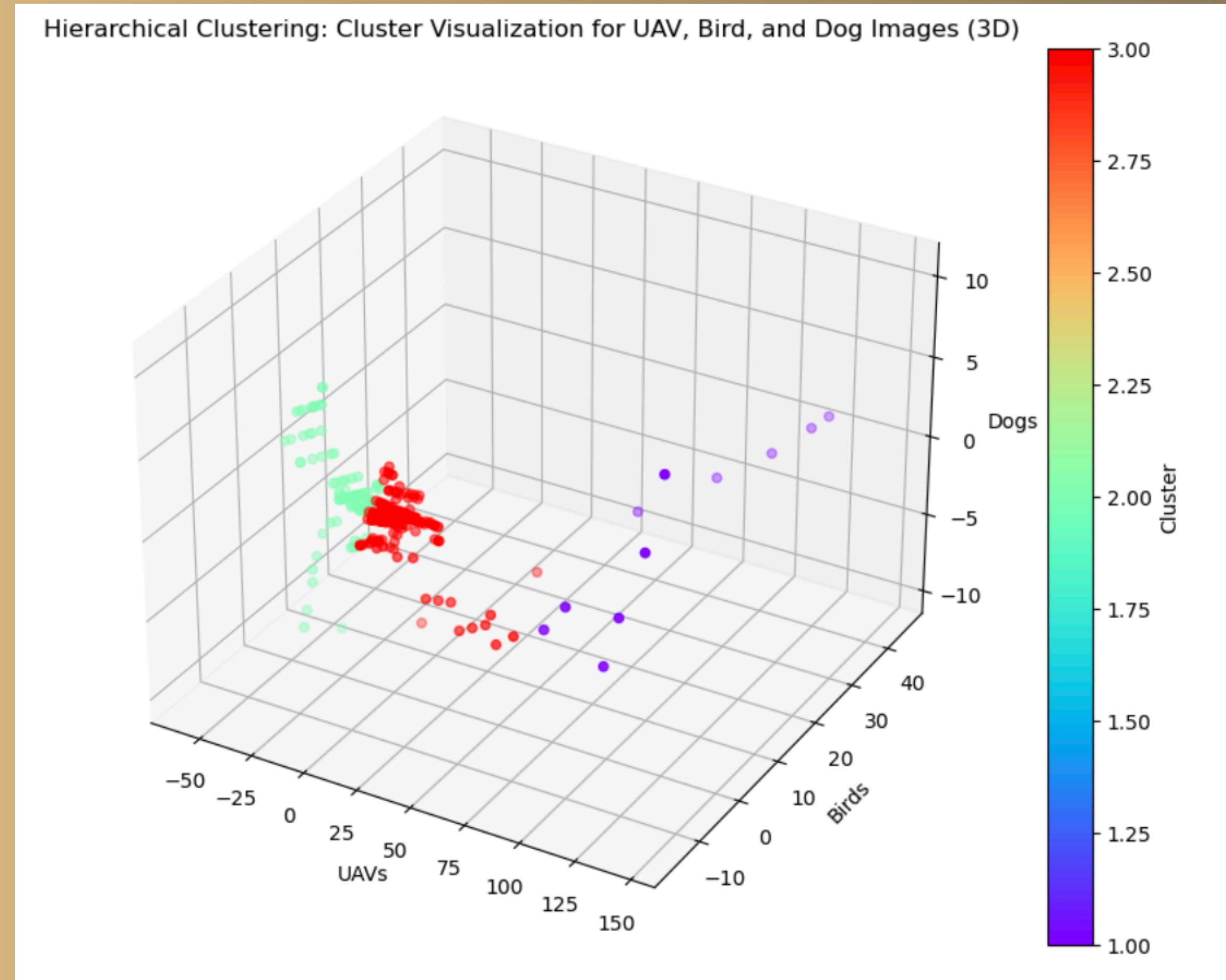
- First unsupervised model approach is a hierarchical clustering model, which creates a cluster for every given feature (image category) that is present.
- The second and third models are anomaly detection models with different base types (NMF and MLP) which are designed to treat UAV data as “anomalies” and group all other image types in a separate classification group.
- The final model is a supervised learning model approach (specifically, a support vector machine classifier) that also performs image classification.

HIERARCHICAL CLUSTERING MODEL



University of Colorado **Boulder**

HIERARCHICAL CLUSTERING MODEL



University of Colorado **Boulder**

ANOMALY DETECTION MODEL - AUTO-ENCODER 01 (MLP)

```
1 # Creating a second anomaly detection model, using MLPRegressor as a base to detect anomalies from
2 # the same dataset (treating UAV image datasets as the anomalies in this case).
3 from sklearn.neural_network import MLPRegressor
4
5 anomalies = uav_preprocessed_data
6
7 # Combining all other datasets (except for UAV data) for training purposes
8 normal_data = np.vstack((bird_preprocessed_data, dog_preprocessed_data))
9
10 # Standardizing the 'normal' data
11 scaler = StandardScaler()
12 normal_data_scaled = scaler.fit_transform(normal_data)
13
14 # Training an autoencoder on the 'normal' data
15 input_dim = normal_data_scaled.shape[1]
16 encoding_dim = 2
17
18 autoencoder = MLPRegressor(hidden_layer_sizes=[encoding_dim], activation='relu', solver='adam', max_iter=500)
19 autoencoder.fit(normal_data_scaled, normal_data_scaled)
20
21 # Standardizing the UAV dataset
22 mlp_anomalous_data_scaled = scaler.transform(anomalies)
23
24 # Retrieving the reconstruction errors for anomalous data
25 mlp_anomalous_data_reconstructed = autoencoder.predict(mlp_anomalous_data_scaled)
26 mlp_reconstruction_errors = np.mean(np.square(mlp_anomalous_data_scaled - mlp_anomalous_data_reconstructed), axis=1)
27
28 # Setting a threshold for anomaly detection
29 threshold = np.percentile(reconstruction_errors, 95)
30
31 # Identifying anomalies in the anomalous dataset
32 detected_anomalies = anomalies[reconstruction_errors > threshold]
33
34 # Calculating the reconstruction errors for anomalous data
35 reconstruction_errors = np.mean(np.square(mlp_anomalous_data_scaled - mlp_anomalous_data_reconstructed), axis=1)
36
37 # Calculate RMSE
38 MLPRegress_rmse = np.sqrt(np.mean(np.square(reconstruction_errors)))
39
40 print("Detected Anomalies:")
41 print(detected_anomalies)
```

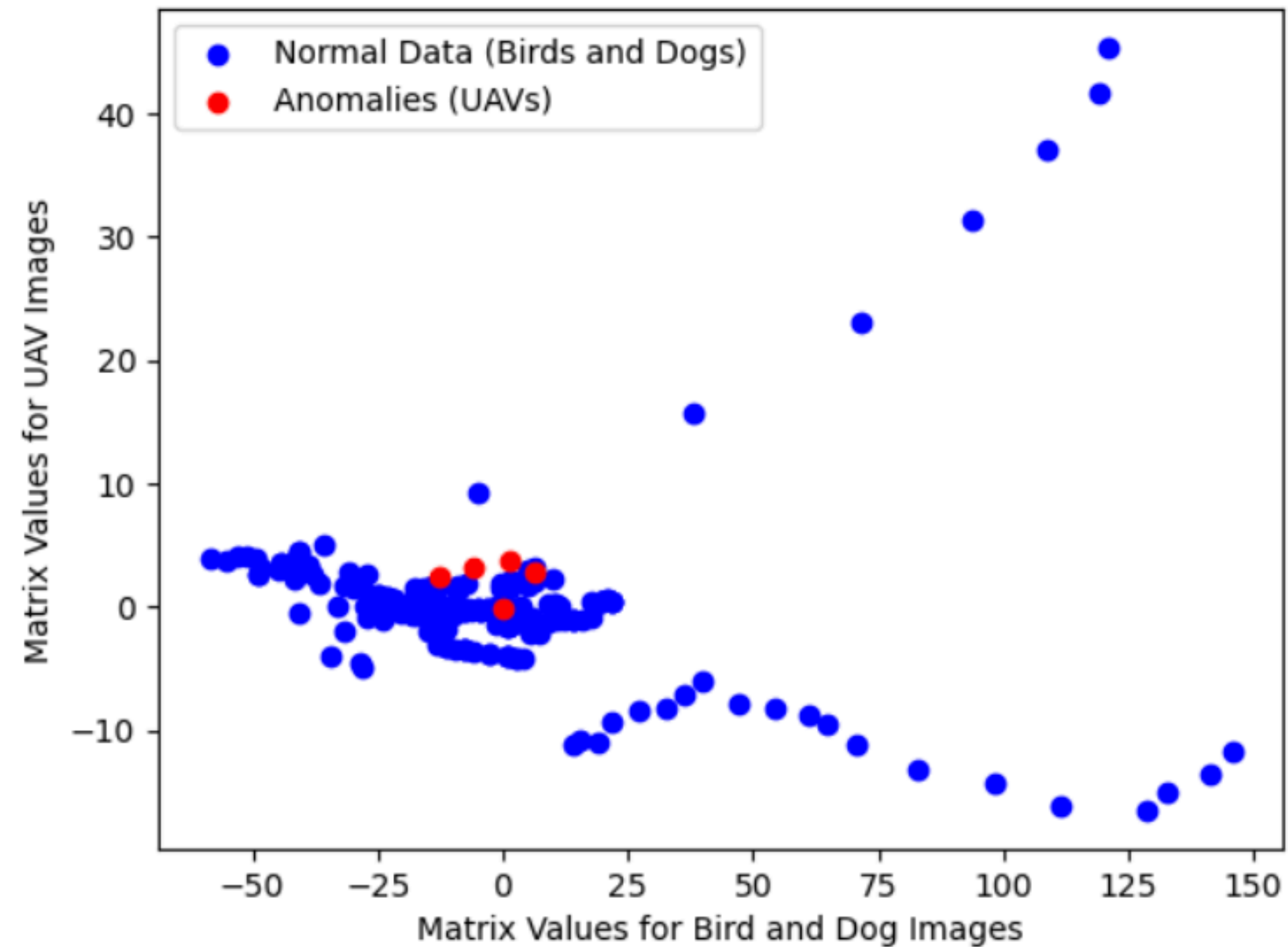
Detected Anomalies:

```
[[-5.72389977  3.1208072  -0.53715342]
 [ 1.17398842  3.76583042 -0.59268238]
 [ 6.47816864  2.84395983 -0.71848694]
 [ 4.7536966   2.68270403 -0.7046047 ]
 [ 3.02922455  2.52144822 -0.69072246]]
```



ANOMALY DETECTION MODEL - AUTO-ENCODER 01 (MLP)

Anomaly Detection using MLPRegressor Autoencoder (Standardized Image Values for Corresponding Image Categories)



ANOMALY DETECTION MODEL - AUTO-ENCODER 02 (NMF)

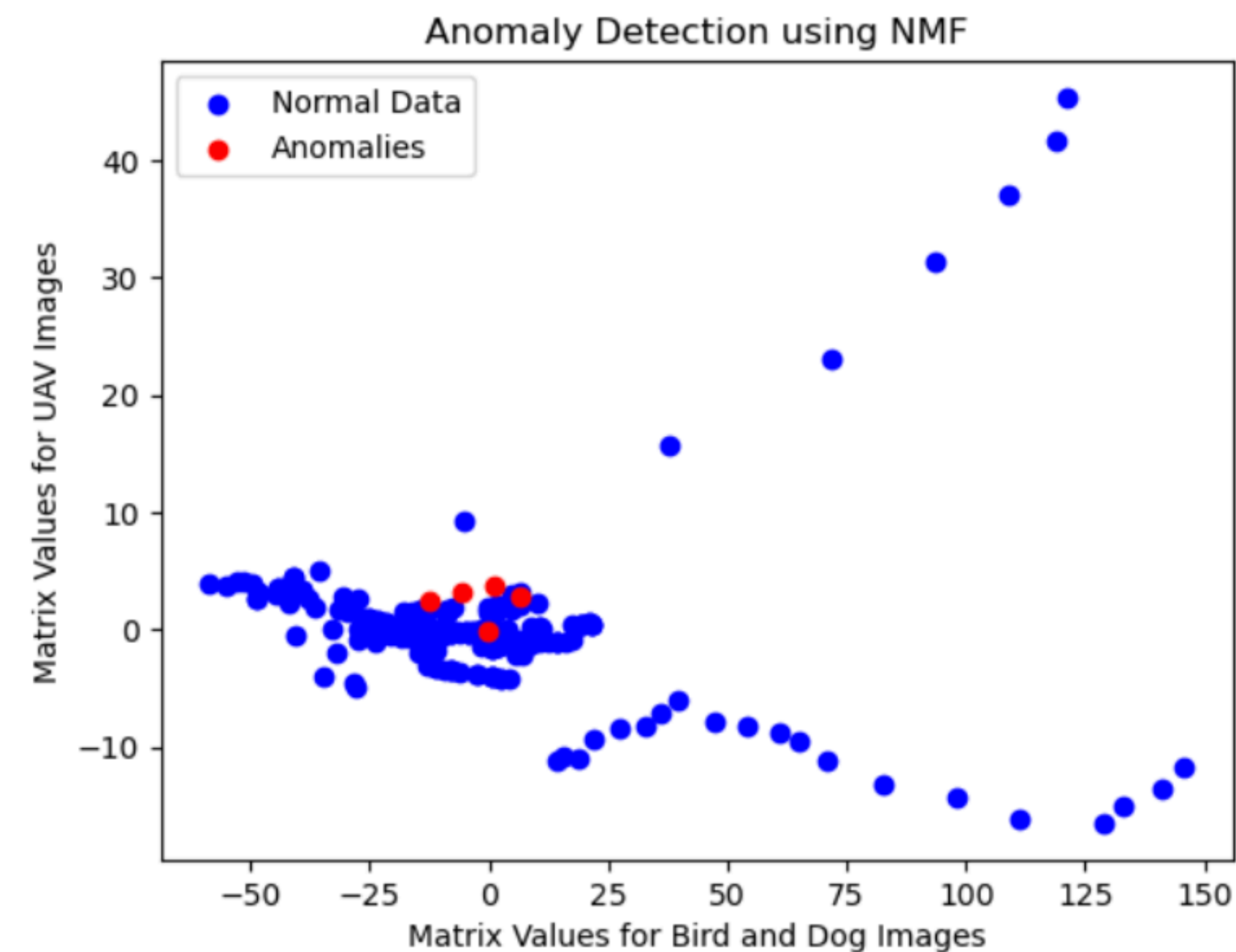
```
# Training an NMF model with the 'normal' data  
n_components = 2 # Number of components for NMF  
  
nmf = NMF(n_components=n_components)  
nmf.fit(normal_nmf_data_shifted)
```



ANOMALY DETECTION MODEL - AUTO-ENCODER 02 (NMF)

Detected Anomalies:

```
[[-12.62178797  2.47578398 -0.48162447]  
 [ -5.72389977  3.1208072  -0.53715342]  
 [  1.17398842  3.76583042 -0.59268238]  
 [  6.47816864  2.84395983 -0.71848694]  
 [-0.17394693 -0.12448924  2.87841013]]
```



University of Colorado **Boulder**

SUPERVISED LEARNING MODEL - SVM CLASSIFIER

```
1 # Creating a SVM classifier model and training it with the same dataset.
2 from sklearn.svm import SVC
3 from sklearn.metrics import accuracy_score
4
5 # Converting the list of lists of matrices into a flattened feature matrix
6 features = [matrix.flatten() for matrix in preprocessed_data[:150]]
7 feature_mat = np.array(features)
8 labels = preprocessed_data[150:300]
9
10 # Splitting the data into training and testing datasets
11 X_train, X_test, y_train, y_test = train_test_split(feature_mat, labels, test_size=0.2, random_state=42)
12
13 y_train = [matrix.flatten() for matrix in y_train]
14 y_train = np.concatenate(y_train)
15 y_train = (y_train > 0.5).astype(int)
16
17 # Training an SVM classifier
18 svm_classifier = SVC(kernel='linear', random_state=5) # Linear kernel for simplicity
19 svm_classifier.fit(X_train, y_train[:len(X_train)])
20
21 # Creating a binary representation for the true and predicted labels.
22 y_test = [matrix.flatten() for matrix in y_test]
23 y_test = np.concatenate(y_test)
24 y_test = (y_test > 0.5).astype(int)
25
26 y_pred = svm_classifier.predict(X_test)
27 y_pred = (y_pred > 0.5).astype(int)
28
29
```



EVALUATION METRICS

- This process involves three main comparisons to determine which is the best quality model for this use case. The first comparison is between the two anomaly detection models, based on their root mean squared error (RMSE) values.
- The second comparison involves qualitative analysis between the best performing anomaly detection model and the hierarchical clustering model. Since the hierarchical clustering model does not involve RMSE metrics, analysis will observe different metrics in the form of within-cluster sum of squares (WCSS) metrics for the hierarchical clustering model.

EVALUATION METRICS (HIERARCHICAL CLUSTERING)

```
1 # Computing the 'within-cluster sum of squares' values for the hierarchical clustering model
2 # over a range of increasing number of clusters to observe the model performance.
3 wcss_values = []
4
5 for num_clusters in range(1, 11):
6     linkage_matrix = linkage(preprocessed_data, method='ward')
7     cluster_assignments = fcluster(linkage_matrix, t=num_clusters, criterion='maxclust')
8
9     cluster_centroids = []
10    for cluster_id in np.unique(cluster_assignments):
11        cluster_points = preprocessed_data[cluster_assignments == cluster_id]
12        cluster_centroid = np.mean(cluster_points, axis=0)
13        cluster_centroids.append(cluster_centroid)
14    cluster_centroids = np.array(cluster_centroids)
15
16    # Calculate pairwise distances between data points and their assigned cluster centroids
17    distances = pairwise_distances(preprocessed_data, cluster_centroids, metric='euclidean')
18    wcss = np.sum(np.min(distances, axis=1)) # Calculate WCSS
19    wcss_values.append(wcss)
20
21 print('Within-cluster Sum of Squares Values for Hierarchical Clustering Model: ', wcss_values)
22
```

Within-cluster Sum of Squares Values for Hierarchical Clustering Model: [5327.626686524979, 4247.639264672412, 3252.270062461683, 2788.043597498734, 2285.323245154893, 2113.812307634568, 1726.6660516382951, 1501.0864706710938, 1426.55325684797, 1280.5615734779371]



EVALUATION METRICS (ANOMALY DETECTION)

```
NMF RMSE Values: [0.5799661121221118, 0.5799661121221118, 0.5799661121221118, 0.5799661121221118, 0.5799661121221118, 0.5799661121221118, 0.5799661121221118, 0.5799661121221118, 0.5799661121221118]
```

```
MLPRegressor RMSE Values: [0.15228291603093527, 0.18179312716755813, 0.10597459151450436, 0.08516477000726493, 0.11487802685919056, 0.051718861525945445, 0.12246666051215643, 0.02257935688840208, 0.04156332181339631, 0.014498409432032688]
```

EVALUATION METRICS (SVM CLASSIFIER)

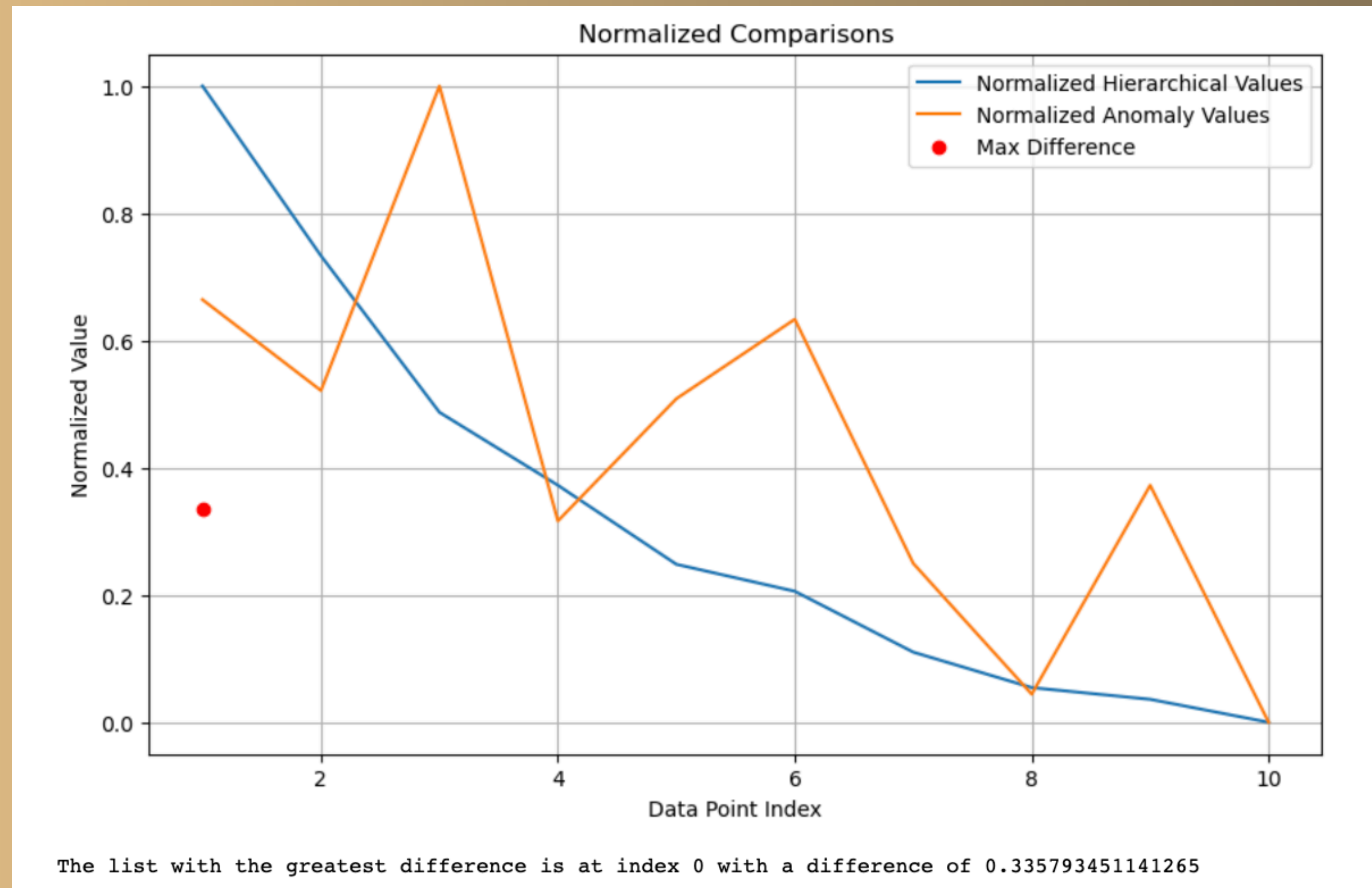
```
Accuracy: 0.7  
Confusion Matrix:  
[[21  0]  
 [ 9  0]]
```

```
SVM Accuracies: [0.6785714285714286, 0.5517241379310345, 0.6896551724137931, 0.6551724137931034, 0.6551724137931034,  
0.7931034482758621, 0.8, 0.6333333333333333, 0.6666666666666666, 0.6333333333333333, 0.7]
```

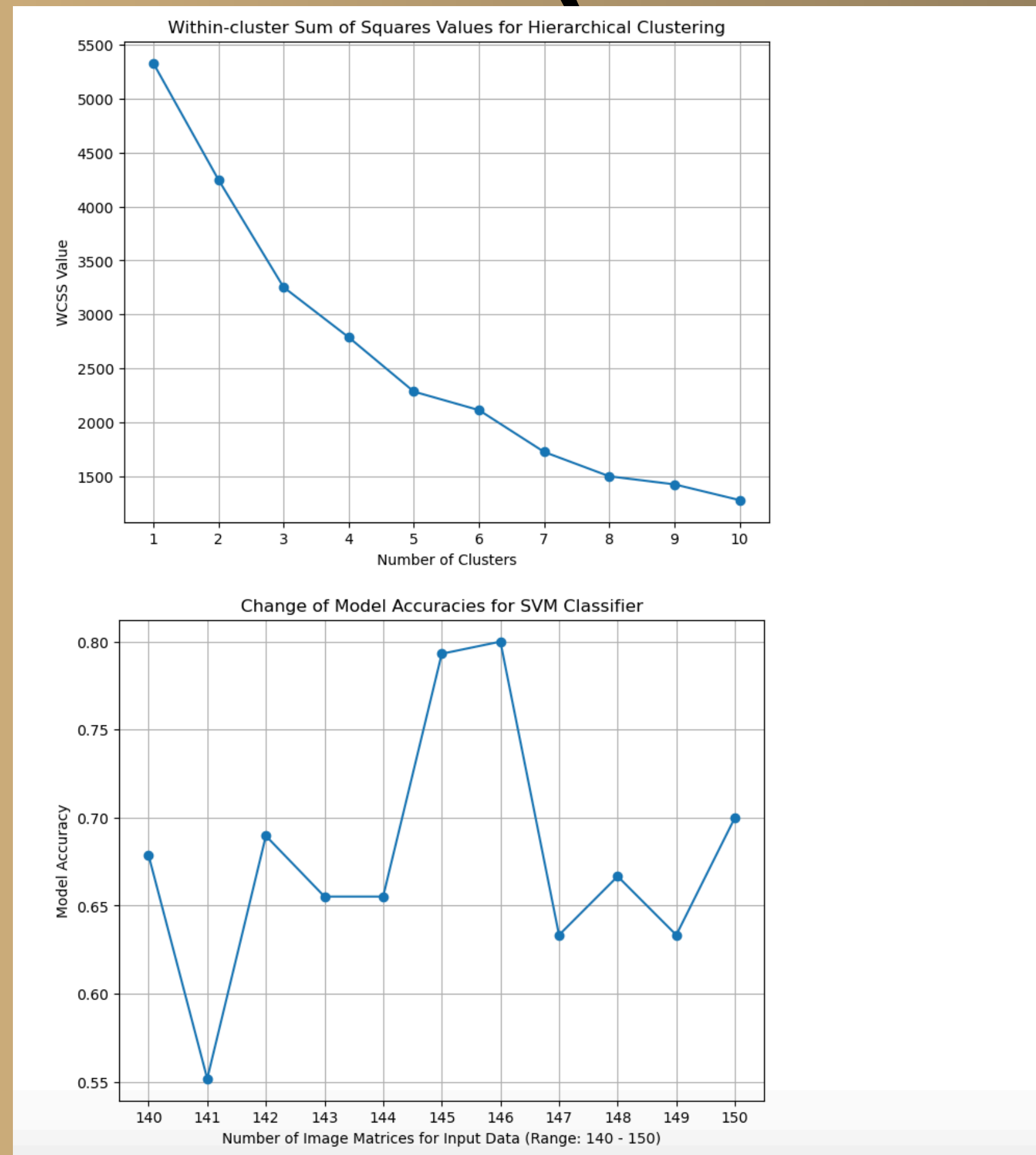


University of Colorado **Boulder**

EVALUATION METRICS (FINAL COMPARISONS)



EVALUATION METRICS (FINAL COMPARISONS)



University of Colorado **Boulder**

RESULTS AND DISCUSSION

- Qualitative analysis was involved in the final comparisons of the models. Even though RMSE metrics could not be generated for the hierarchical clustering model, there was still enough information to receive a general perception on the overall behavior and performance of each model.
- The visualizations indicate that there are inconsistencies in the error metrics and accuracies when training the following models through multiple iterations: both anomaly detection models and the SVM classifier model. In regards to the hierarchical clustering model, this was the only model that showed consistency and gradual improvements across multiple iterations.

LESSONS LEARNED

- Data understanding and the model selection process were very crucial aspects of this project.
- Testing the same machine learning models with larger scales of data can determine robustness of the given models (but could also impact overall performance and runtime).
- Making use of different types of evaluation metrics



CONCLUSION

- In conclusion, the hierarchical clustering model was found to be the best model in terms of consistency and overall performance.
- Further tests and model optimization can improve the performance of all the other models involved.
- Incorporating other types of clustering models (such as K-means clustering) could yield different but better results in the search for best models to use for UAV image classification

THANK YOU



University of Colorado **Boulder**