

# Austin Animal Shelter

Predicting Shelter Outcomes  
for Cats and Dogs

Madeleine Sharp | DSC630 | Final Project Paper

---

## Introduction

Each year, approximately 6.3 million companion animals enter U.S. animal shelters nationwide (ASPCA, 2022). Of those 6.3 million shelter animals, approximately 920,000 are euthanized annually (although, this number has declined within the past eleven years). While myriad shelters are not no-kill (and therefore, do euthanize animals past a certain time point), there are shelters in existence that are no-kill shelters and actively seek to try to do all they possibly can for an animal outside of euthanizing the animal.

One of these shelters is the Austin Animal Center (AAC) in Austin, TX. This shelter provides care and shelter to more than 18,000 animals annually and is also heavily involved in various activities and initiatives for the protection and care of abandoned, at-risk, or surrendered animals. The Austin Animal Center works with county, city, and state-wide parties on these efforts in an endeavor to better the lives and outcomes for those animals in need of care and a home.

Given the large number of animals residing in shelters, alongside the large number of those animals euthanized, finding rehoming solutions for animals and shelter pets in need is a critical and key issue for animal shelters across the nation. Even more so for those shelters that are no-kill shelters, given that they keep animals until they can be moved on to that next step of their lives, whether that be adoption, rehoming, shelter transfer, etc. Of course, and unfortunately, even at no-kill shelters, certain instances may exist in which an animal is euthanized due to being extremely sick or in extreme pain. Generally, however, this method is only proceeded with once all other potential methods have been considered or exhausted.

Solving the animal shelter adoption problem is imperative, and those who would most likely benefit from such an endeavor include the animals, people, facilities, structures, and institutions in the animal care sector/industry. Keeping pets off the streets and in good homes is imperative for the health and well-being of these animals. Additionally, their health and well-being possesses a ripple effect outwards towards all communities and institutions that serve animals or that may be impacted. These communities and institutions include, but are not limited to, veterinary offices, animal control offices, animal shelters, public health offices, etc.

One of the ways in which the power of data science and machine learning can help to solve the animal shelter adoption problem is via a predictive approach. More specifically, predicting animal outcomes based on an animal's features to determine (ahead of time) the future of that animal. This would be helpful for planning purposes for shelter facilities, as well as determining resource needs and allocations. In this vein, my project focus is centered on prediction in an endeavor to help AAC determine the future outcomes of their shelter animals.

Overall, the scope of my project outlined within this paper consisted of utilizing the AAC data and its characteristic/demographic variables to ascertain "next step" outcomes of the animals. More specifically, outcome type is a variable within the dataset, and each animal outcome is a potential classification (adopted, transferred, etc.) of that. Therefore, the predictive focus of my project was classification, given this is a classification problem. I wanted to know, based on features from the dataset, if the outcome of an animal (post-shelter) could be predicted.

Via the use of AAC's data, I conducted an exploratory data analysis (EDA), and utilized classification machine learning (ML) methods to predict animal outcomes at AAC's shelter.

## Methods

### *Data, Data Dictionary, and Variables*

As a part of Austin, Texas's Open Data Initiative efforts (a city effort), the AAC has made its data available for public use (City of Austin, TX, 2022). To obtain the data, there are a few components: the data itself (collected and managed by the AAC), extracting that data (via the Open Data Portal's API and the use of a Socrata script, which is what the API is powered by) (Socrata, 2022), and then structuring that data in such a way within a Python pandas data frame for usage.

The various tools and resources I utilized to obtain the data for this project include (please reference the links below):

- <https://data.austintexas.gov/>
- <https://data.austintexas.gov/Health-and-Community-Services/Austin-Animal-Center-Outcomes/9t4d-g238>
- <https://gist.github.com/aschleg/54bf7ed55c2383f3ba1f338b8116a77b>
- <https://data.austintexas.gov/resource/hcup-htgu.json>
- SODA - the Socrata Open Data API (SODA), which provides programmatic access to the dataset for use (<https://dev.socrata.com/consumers/getting-started.html>).

This specific dataset's variables included:

- animal\_id - This refers to the animal's specific ID number.
- datetime - Date and time of when the animal was admitted to the shelter.
- monthyear - Month and year of when the animal was admitted to the shelter.
- date\_of\_birth - Date of birth of the animal, if known.
- outcome\_type - Outcome type, or classification.
- outcome\_subtype - Outcome subtype (if exists), or classification.
- animal\_type - Type of animal (dog, cat, horse, etc.)
- sex\_upon\_outcome - Sex of the animal at the time of their "outcome."
- age\_upon\_outcome - Age of the animal at the time of their "outcome."
- breed - Animal's breed.
- color - Color of animal's fur/coat.
- name - Animal's name.

This data spans back to the year 2013 and up to the present date. It is a real-time and streaming dataset.

### *Obtaining and Extracting the Data*

The data for this project was obtained via the use of the SODA API script (the API used to access the AAC data is powered by Socrata). The API possessed a call limit of 1,000, which meant that the results needed to be garnered 1,000 results at a time (this is also noted within the API documentation on the website). At the time of my completion of this step, the documentation on the Austin, TX AAC website page indicated that approximately 140K total rows/records of data existed within the dataset. To obtain the data, an endpoint for where the data is hosted required creation as defined in the API documentation on the website. I needed to define the number records to call (1,000) as well as the number of pages to "flip through" for getting the data. For pages, per the website, setting this "paging" threshold at 200 was the most beneficial. Once the data was extracted, I read it in as a pandas data frame, and then also saved it as a .csv file for potential future use and retrieval.

### *Cleaning and Transforming the Data*

Following data extraction, the data was required to undergo some changes in order to be the most useful for analysis and machine learning model building. The steps taken to clean the data included the following:

- I assessed for any duplicate entries of the animal\_id variable/
  - I determined if each animal\_id and datetime combination was a unique combination (which would signify a separate intake for the same animal).
- I assessed for unique values within each variable in the dataset to ensure legitimacy.
- I assessed for any missing variables and handled those via removal – since not many missing values existed (less than 50 in a dataset with approximately 140,000 entries), I elected to remove all missing values and their respective rows/entries.
- I converted date and time related variables to a usable format (datetime).
- I split the data into separate data frames: one for cats and one for dogs.
- I completed some feature engineering so that the variables were in the best shape for usage during the model building and running phase.
  - This is discussed more in detail in the Feature Engineering section of this paper.

### *Exploratory Data Analysis (EDA)*

For the EDA portion of this project, I completed two phases of EDA. The first phase was an overall EDA for the dataset as a whole, prior to any feature engineering or identification/selection. The second EDA phase occurred later in my project, following additional feature work, and included a separate EDA for cats and dogs.

For the first EDA, I mainly assessed variable distributions, with the variables of focus being animal type, outcome type, sex of the animal at the time of outcome, and animal number of stays at the clinic.

Secondly, following additional feature variable work, I focused in on cats and dogs EDA. While I used the same EDA techniques for both, and largely assessed the same variables, I did so separately – cats first, and then dogs. For these EDA portions, the variables of focus included assessing outcome type and outcome subtype by age of the animal as well as sex of the animal, and outcome types across time (from 2013 to present date).

The results of the EDA are discussed in the Results and Evaluation section of this paper, and the visualizations are included in the Appendix at the end of this document, beginning on **page 10**. Each EDA section (overall, cats, and dogs) is labeled respectively.

### *Feature Engineering*

To use this data to the greatest capacity, and for the best machine learning model outcomes, I elected to do some feature engineering. This allowed me to obtain the most from the features in the data set and allowed me to engineer them in such a way that they could be the most useful and “read in” more easily by the machine learning model.

The specific variables that were transformed for feature engineering (for both cats and dogs) included:

- sex\_upon\_outcome (two new features created: sex and Spay/Neuter).
  - This variable was essentially separated into two new features.
- age\_upon\_outcome (Periods, Period Range, outcome\_age\_(days), outcome\_age\_(years), Cat/Kitten (outcome) & Dog/Puppy (outcome), sex\_age\_outcome, age\_group, dob\_year, dob\_month, and dob\_monthyear were engineered largely from this variable).
  - This variable was essentially split out for the animal’s age and its respective shelter-handling dates. This allowed for easier interpretation for analysis.
- breed (multiple new features created: cfa\_breed, breed1 (breed\_1), breed2 (breed\_2), and domestic\_breed).
  - These were engineered by using a list of cat/dog breeds from The Cat Fanciers’ Association & The American Kennel Club and then splitting the breed variable out into a main breed plus a secondary breed and whether the animal was a domestic breed or not.
- color (here, I engineered coat\_pattern, an updated color variable, color1(color\_1), and color2 (color\_2)).
  - This broke up the animal’s color variable into color and coat pattern (so two new variables split from the original one). A second color variable was also used if it was present.

Overall, the original variables were kept where necessary, and dropped where unnecessary. This was largely dependent on the feature engineering. If the feature engineering satisfied the goal in and of itself using the original feature, then that original feature was dropped. If the original feature remained to serve an important purpose apart from the newly engineered features, it was kept.

### *Feature Identification and Selection*

In addition to completing some feature engineering work, I elected to also implement some feature identification and selection methods. This involved using a random forest classifier to identify the top fifteen (15) features, and Principal Component Analysis (PCA) for feature selection.

With the random forest classifier, for cats and dogs, different features appeared to be most important. For cats, these were:

- Spayed females.
- Neutered males.
- If it is a cat.
- If it is a kitten.
- Intact female.
- Intact male.
- The outcome hour being 0 (midnight) (interesting).
- Having a tabby coat pattern.
- Being born in April (interesting).
- Being brown in color.
- Unknown sex (interesting).
- Being approximately 2 years of age or under.
- Outcome month being September (interesting).
- Outcome hour being 13 (1PM) (interesting).

For dogs, these were:

- Spayed females.
- Intact males (interesting).
- Neutered males.
- Intact females.
- If the breed is a pit bull.
- Under 3 years of age.
- If the dog is an adult dog.
- If the dog is a puppy.
- If the outcome month is June (6 - summer).
- If the breed is a Labrador retriever.
- If the breed is a chihuahua.
- If the dog is between the ages of 2 and 5.
- If the dog was born in April (4).
- If the outcome hour was 0 (interesting).
- If the outcome hour was 12.

Next, I implemented Principal Component Analysis (PCA). PCA is a technique that is commonly used in supervised ML algorithms across a variety of applications. It is primarily used for dimensionality reduction. Dimensionality reduction reduces the number of input variables (features) – therefore is a form of feature selection via feature reduction. I completed PCA both

for the cat and the dog model building. The goal of using PCA in these instances was to allow for dimensionality reduction to ensure the features to be used would be the most robust and the best fit for feeding into the model.

Overall, the random forest classifier helped to identify important features in the dataset, and PCA helped to reduce those via dimensionality reduction.

### *Machine Learning and Model Building*

Given that the dataset I used involved data types that would require an aspect of classification for predictive measures, I sought to use a classification machine learning model for predicting post-shelter animal outcomes. To complete this task, I built two models: one for cats and one for dogs. I determined that by creating and running a separate model for each animal group, I would obtain the most accurate results for that respective group. The following portion of this section describes my process for the ML model building and this process applies to both the cat and dog models.

With respect to the model classification groups and predicting the target (outcome type), I needed to encode the data via the use of dummy variables for the features (all other variables) in the dataset to be used as potential predictors. I then factorized the target variable, and next split the data into the training and testing sets. I also used a standard scaler to scale the data. Following that, I used the random forest classifier mentioned in the Feature Identification and Selection section of this paper above. Then, coupled with my PCA, I created a feature union in which I used PCA and Select K Best to obtain and pull together the best features for use within the model. Lastly, I moved onto the building of my ML model.

Based upon my research regarding the best type of model to use for this data, as well as the best parameters, I selected to build a pipeline that included my feature union and a random forest classifier. Within the random forest classifier, the parameters consisted of the below:

- `n_estimators=1000, criterion='gini', class_weight='balanced_subsample', bootstrap=True, oob_score=True`).

Again, the determination for these came based upon research regarding other classification model case examples, and some that even focused on animal shelter classification (Kaggle, 2016; Vickery, 2019).

After composing and compiling the pipeline model with its parameters, the model was fit to the training data. I used a cross validation Stratified K Fold of 2, and then printed the results of the cross-validation score and prediction score. Next, I used the model for prediction on the test data and then printed the accuracy score of the model. Lastly, I printed a normalized and non-normalized confusion matrix to view the accuracy results of each of the potential outcome types. I also printed a classification report.

For both cats and dogs, I did this twice: once with all outcome types, and once with two main outcome types – Adoption and Transfer. The secondary models, for both cats and dogs, performed better than the first by ten percentage points. The results summaries and information for each of these models can be seen below in the Results and Evaluation section, as well as in the Appendix of this document, beginning on **page 17**. Each section is labeled respectively.

## Results and Evaluation

### EDA

From the EDA, I was able to obtain insights for the dataset distributions overall, as well as insights regarding the specific cat and dog data within the data set.

In the initial EDA, the results indicated that the majority of the dataset was comprised of dogs, the main outcome type was a transfer (shelter transfer), most of the animals were neutered males, and the majority of the animals had only visited the shelter once (so not many repeat visits or intakes).

With respect to the cat-specific EDA, I was able to drill-down a bit more within the data. Overall, the cat EDA indicated the following findings:

- Intact males and females of kitten age (six months or less) were transferred more than any other outcome.
- The results were consistent with general shelter practice as intact cats are generally spayed or neutered before becoming available for adoption.
- Kittens were more likely to be transferred to partner facilities, perhaps for receiving treatments (spayed/neutered, vaccines, etc.).
- Neutered/spayed kittens were adopted more frequently than adults.
- Kittens were typically adopted faster (less time in the shelter) than adult cats.
- The cats entering the shelter were usually under five years of age.
- Adoptions and transfers tended to be seasonal, peaking in the summer months (mid-year) and then dropping back down until the subsequent summer season.
- Other outcomes such as return to owner were less seasonal and not as easy to track/glean meaning from across time.
- There did not appear to be any significant relation in the outcome to the age group of the cat, other than if the cat was of kitten age or not.

From the dog-specific EDA, the drilled-down findings indicated that:

- Dogs under two years old made up the initial large majority of records.
- There were approximately the same number of female and male dogs within the dataset.
- Dogs under five years of age secondarily made up the majority of records.
- The most common outcomes were adoption or transfer to a partner facility.
- Euthanasia outcomes greatly decreased overall across time since pre-2014.
- Lots of deaths, yet not many missing or return to owner outcomes.
- Neutered males and spayed females appeared to largely comprise the bulk of the adopted outcome.
- Generally, it looked like most of the dogs were approximately age 15 and under.
- Neutered and spayed dogs were most often transferred.
- Dogs were adopted more frequently than puppies.
- Most dogs and puppies were transferred to partner facilities.

Please note that the visualizations for these findings are included in the Appendix, at the end of this document, beginning on **page 10**.



## *Machine Learning Model*

Overall, the results and model outcomes for this project were relatively strong. While in some instances, they were not as strong as I had hoped, given the amount of data in use for this project, the performance of the ML models was still significant.

Firstly, I created the two cat ML models, as described in the Methods section of this paper. The first model focused on all the outcome types, which included Transfer, Adoption, Euthanasia, Died, Return to Owner, and Missing. The model, when set at a 70/30 train and test split, performed with approximately 81% accuracy. In addition to this model, I also built and ran a model that focused on the two main outcome types, which were Adoption and Transfer. This model had increased and improved performance, with an accuracy score of about 91% (recall was also strong). To view the full results from the confusion matrices and the classification reports, please refer to the Appendix at the end of this document, beginning on **page 17**.

Secondly, I also created two dog ML models (also as described in the Methods section of this document). The first model, like the cat model, focused on all the outcome types. This model, when set to a 70/30 train and test split, performed with approximately a 71% accuracy. The secondary model, focusing only on Adoption and Transfer outcome types, performed at about 81% accuracy (recall was strongest here as well). Again, to view the full results, please refer to the Appendix at the end of this document, beginning on **page 20**.

Overall, the secondary ML models for both cats and dogs performed significantly better than the first models that included all the outcome types. The cat model performed better than the dog model; this could be due to the differences in the number of and potential unique value types for each of the feature variables for cats and dogs, respectively.

## **Conclusion & Final Takeaways**

In conclusion, these predictive machine learning models indicated that they could be useful for predicting animal shelter outcomes before those outcomes might occur. This is helpful for the Austin Animal Center (and likely other shelters as well) for planning and decision-making purposes. The models performed quite well, especially as it pertains to real-world applications, and they are a good foundation for improved data considerations and decisions down the line.

One area in which this could have been made to be even more robust is by splitting out the data to account for “past” and “future.” In this project, I did not split the data according to date to test the models. Since the data is streaming, an additional step moving forward could be to split the data at a certain time point, and then see how well the model did relative to the actual data used for the testing portion/split.

In any case, the EDA granted insights into the status of the various shelter and animal demographics, and the machine learning models helped to predict the specific outcomes for cats and for dogs at the AAC.

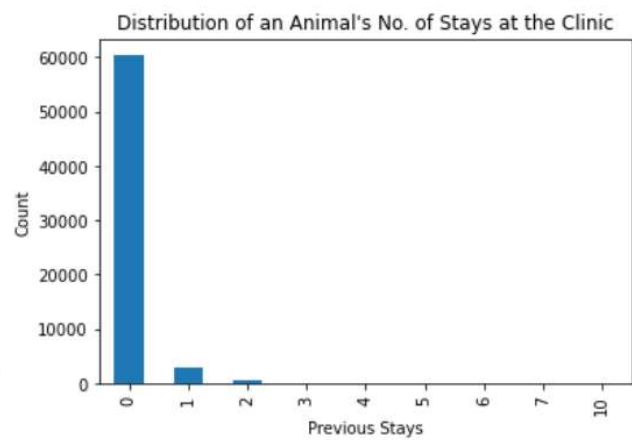
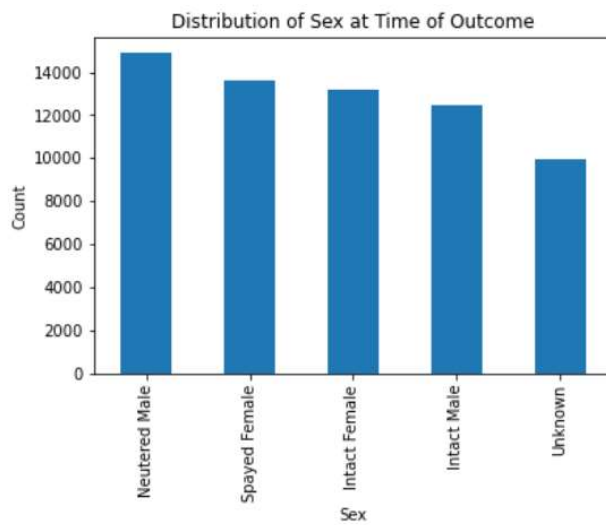
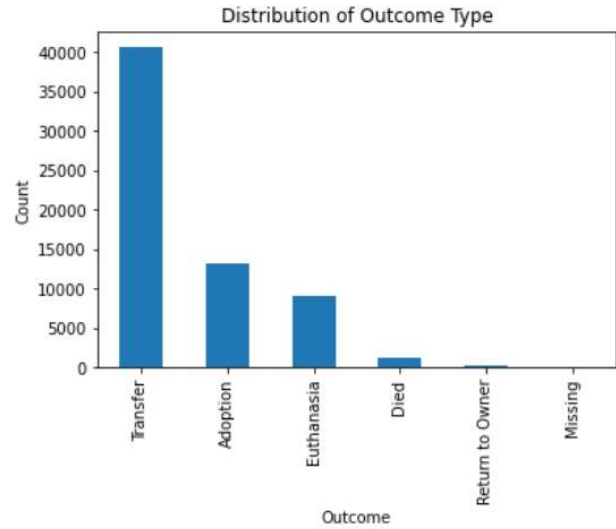
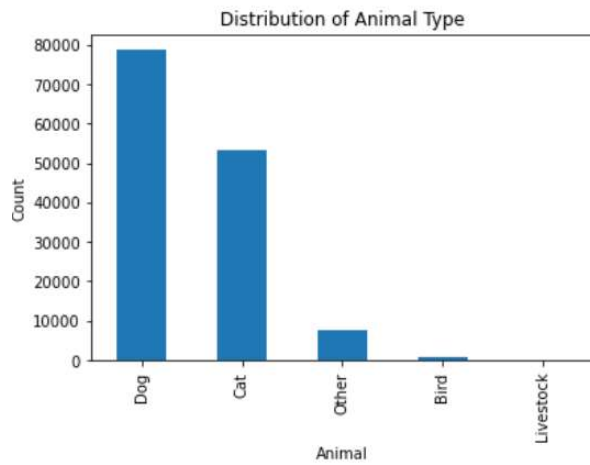
### References:

- ASPCA. (2022). *Pet statistics*. ASPCA. <https://www.aspca.org/helping-people-pets/shelter-intake-and-surrender/pet-statistics>
- City of Austin. (2022). *Austin Animal Center outcomes*. City of Austin.  
<https://data.austintexas.gov/Health-and-Community-Services/Austin-Animal-Center-Outcomes/9t4d-g238>
- Kaggle. (2016). *Shelter animal outcomes: Help improve outcomes for shelter animals*. Kaggle.  
<https://www.kaggle.com/c/shelter-animal-outcomes>
- Schlegel, A. (n.d.). Simple script for extracting Socrata Open Data Access (SODA) datasets.  
GitHub Gist. <https://gist.github.com/aschleg/54bf7ed55c2383f3ba1f338b8116a77b>
- Socrata Developers. (2022). *Consumers: Getting started*.  
<https://dev.socrata.com/consumers/getting-started.html>
- Vickery, R. (2019). *Predicting animal shelter outcomes*. Medium.  
<https://medium.com/vickdata/predicting-animal-shelter-outcomes-4c5fad5dbb4f>

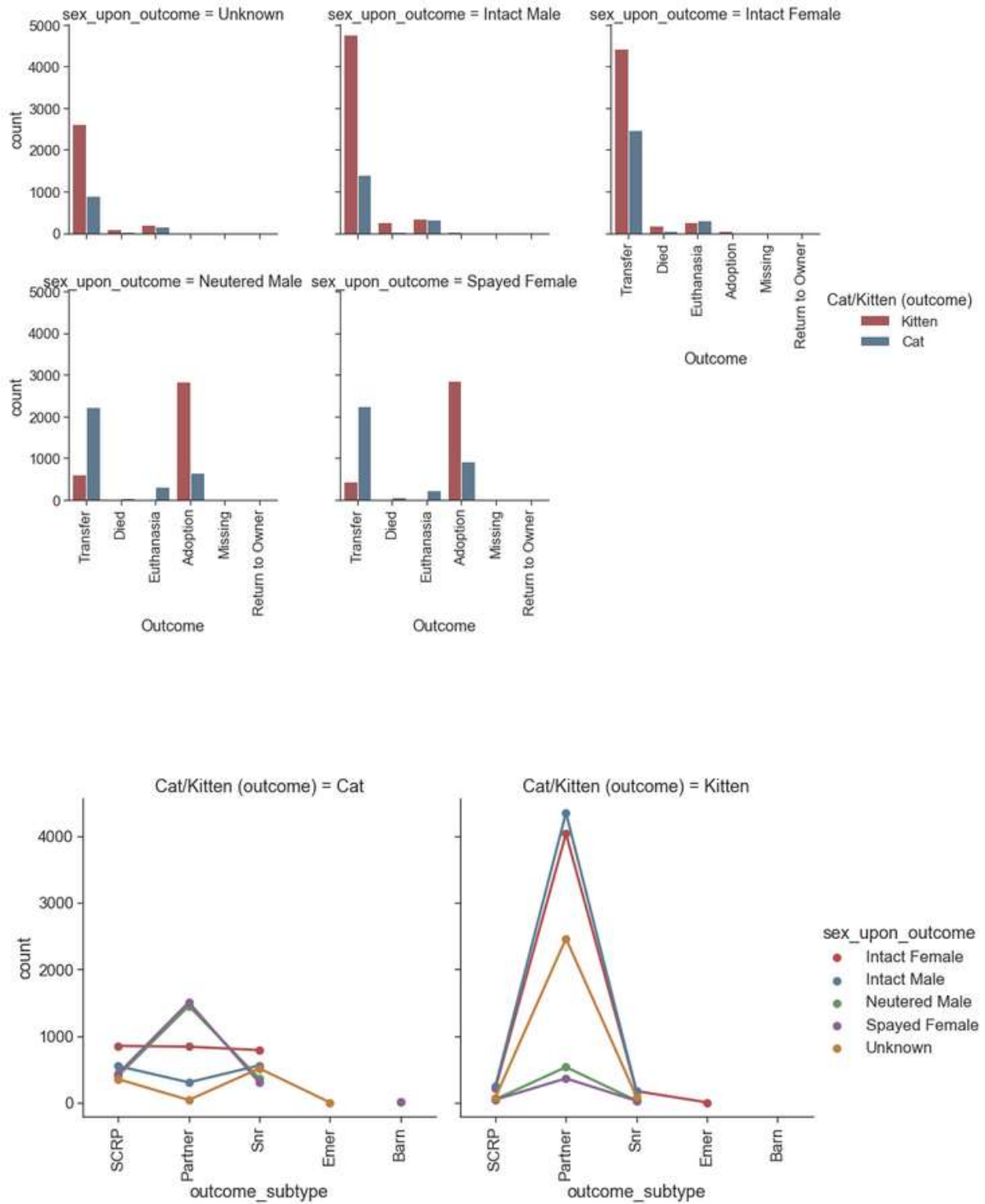
## Appendix

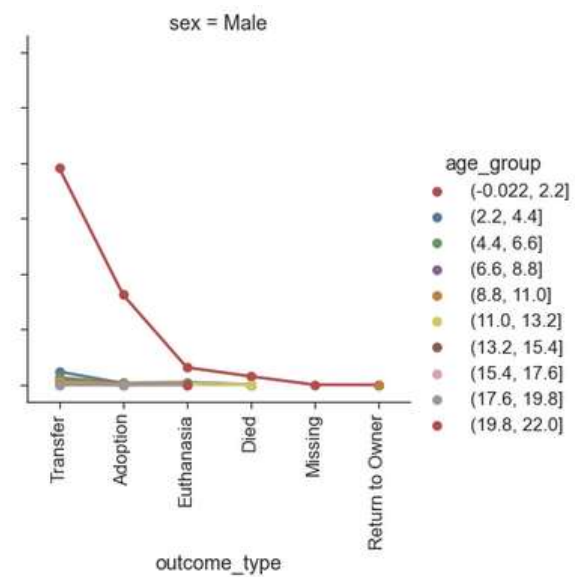
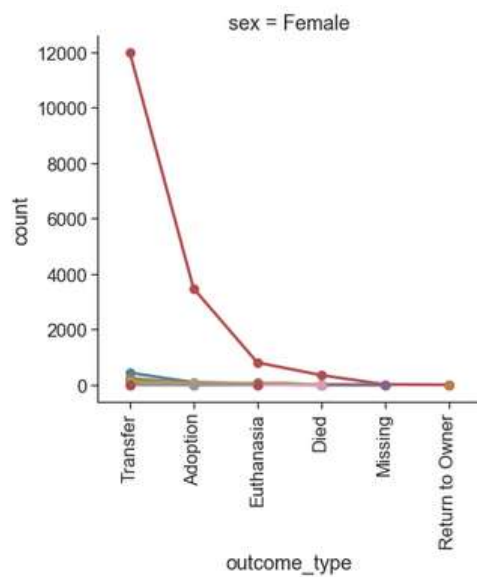
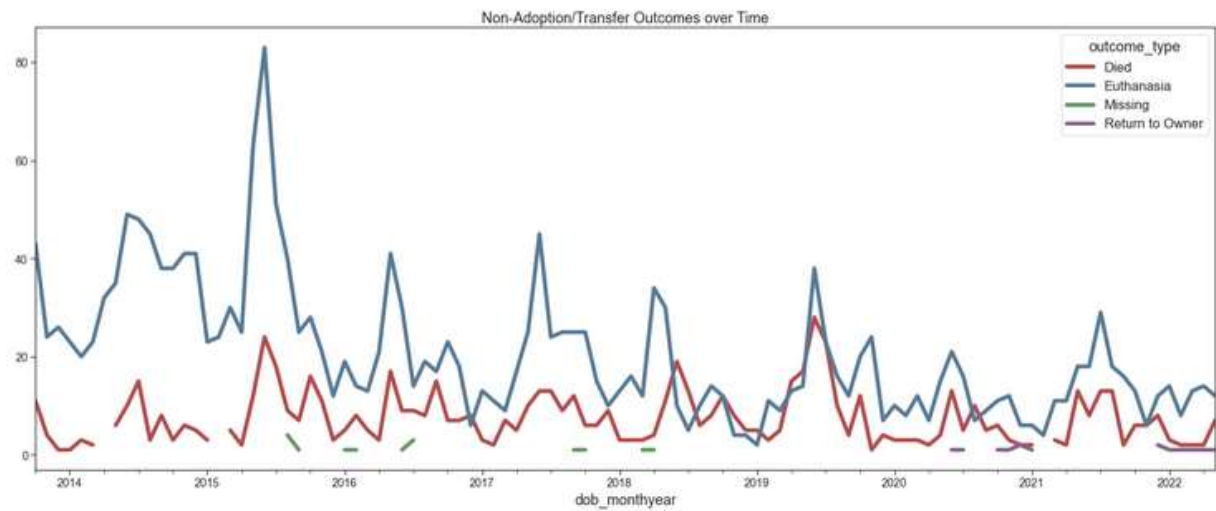
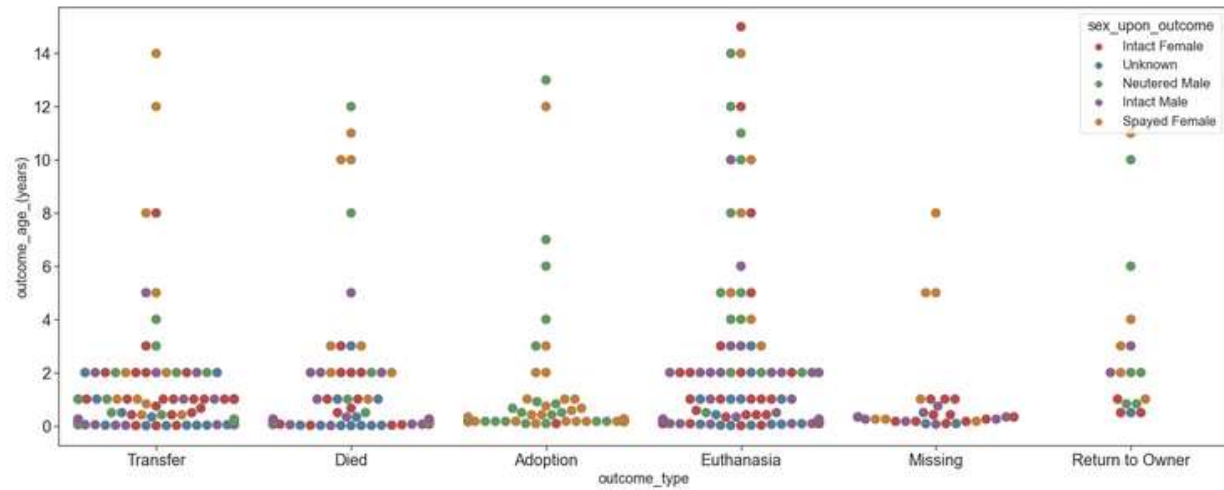
### EDA Visualizations

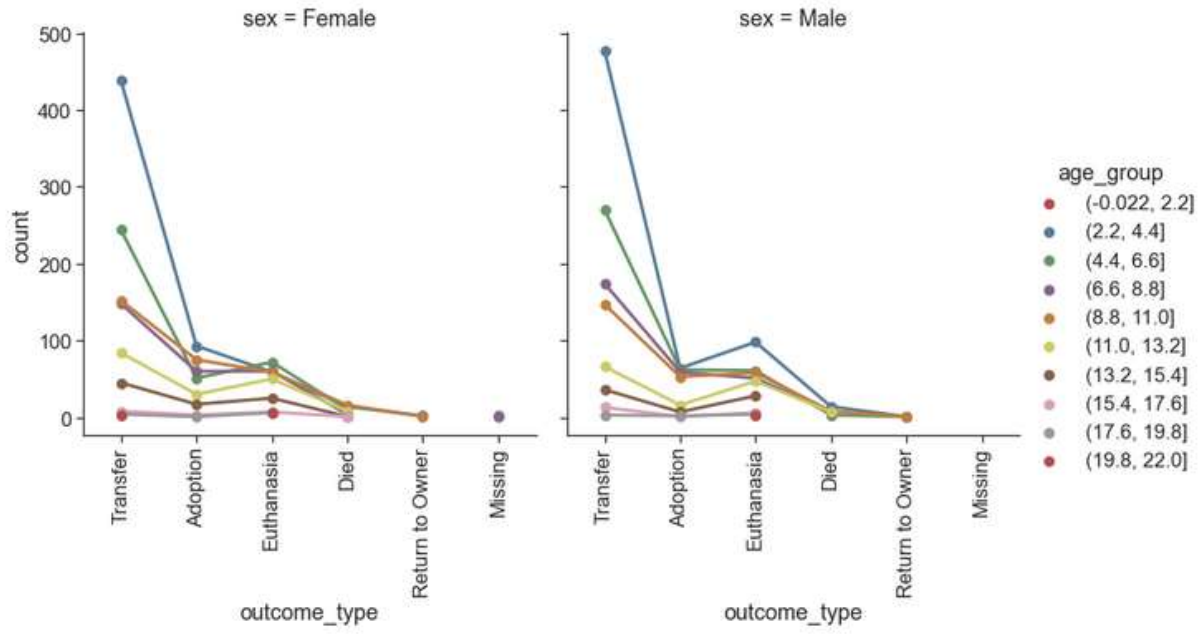
#### Overall EDA:



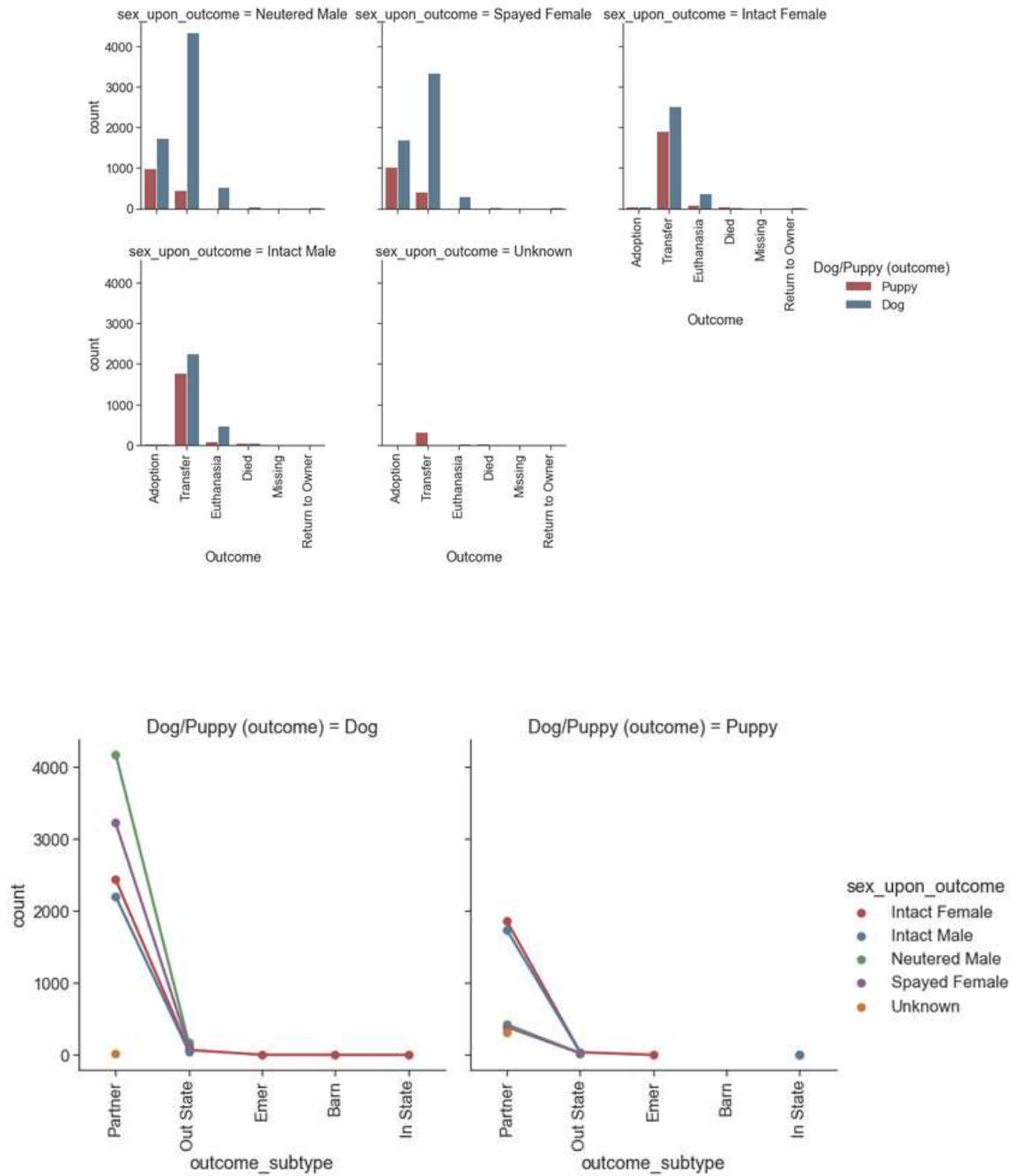
## Cat EDA:

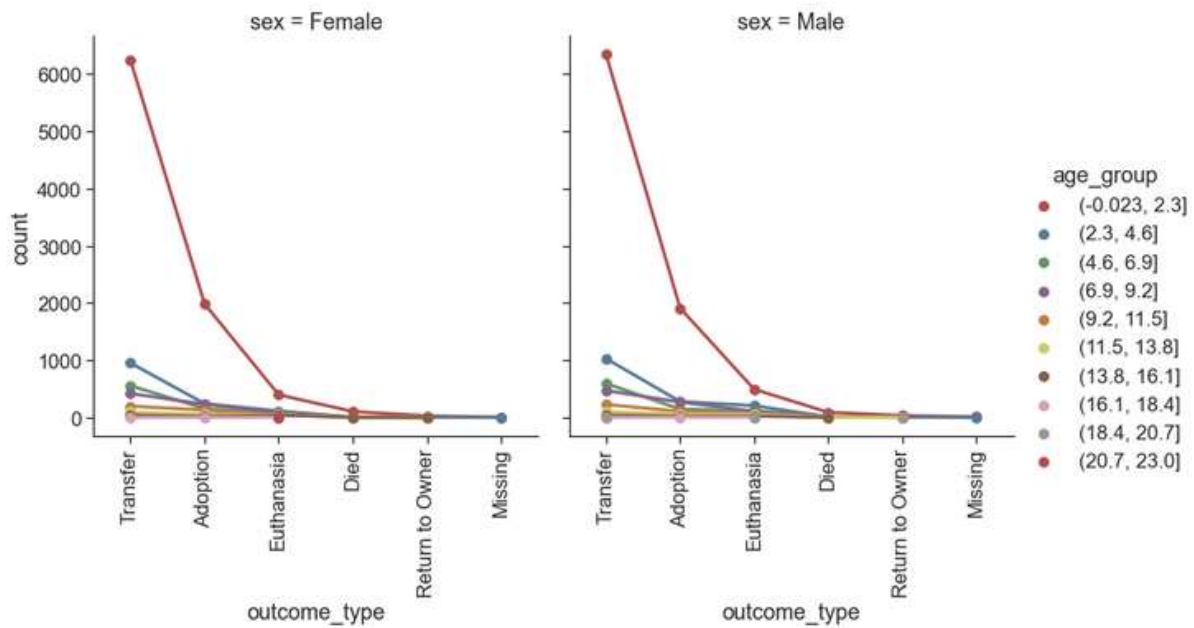
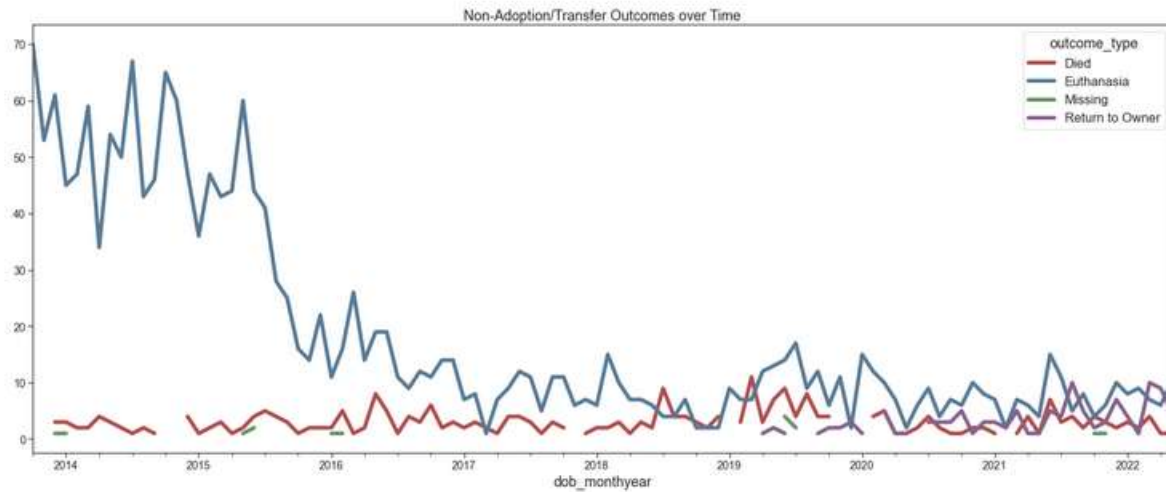
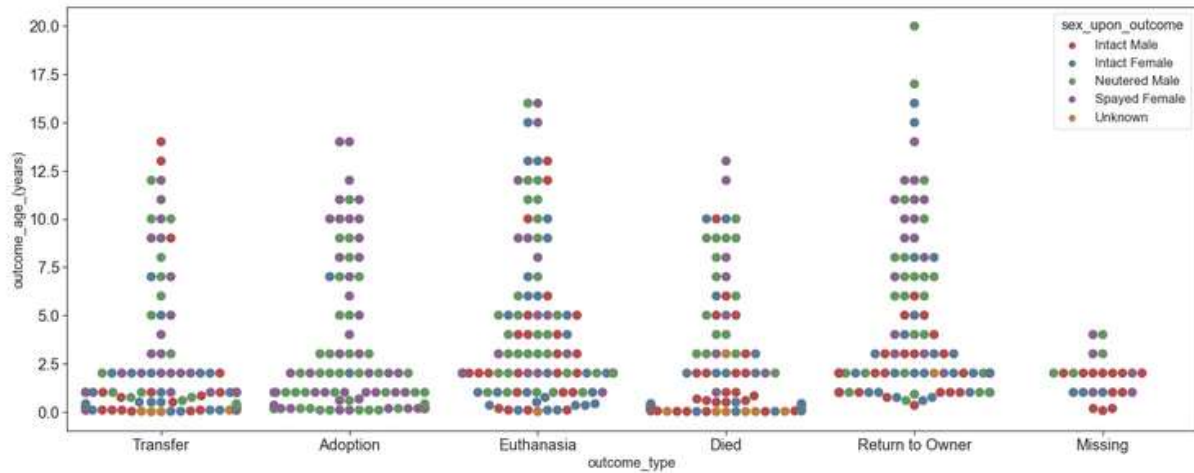




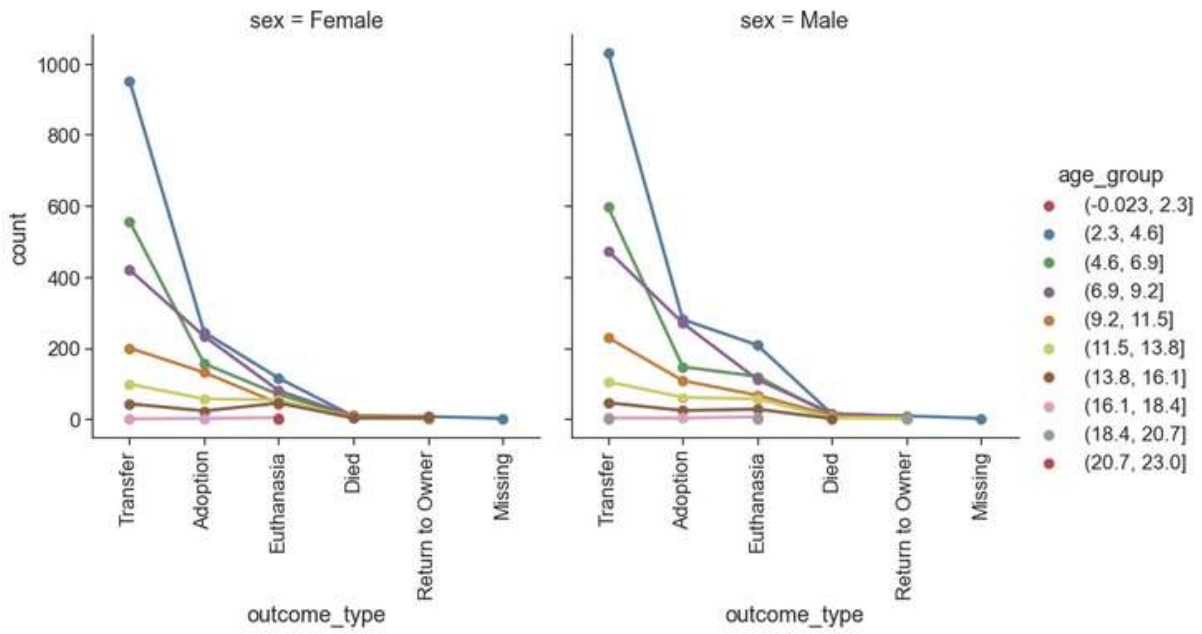


## Dog EDA:



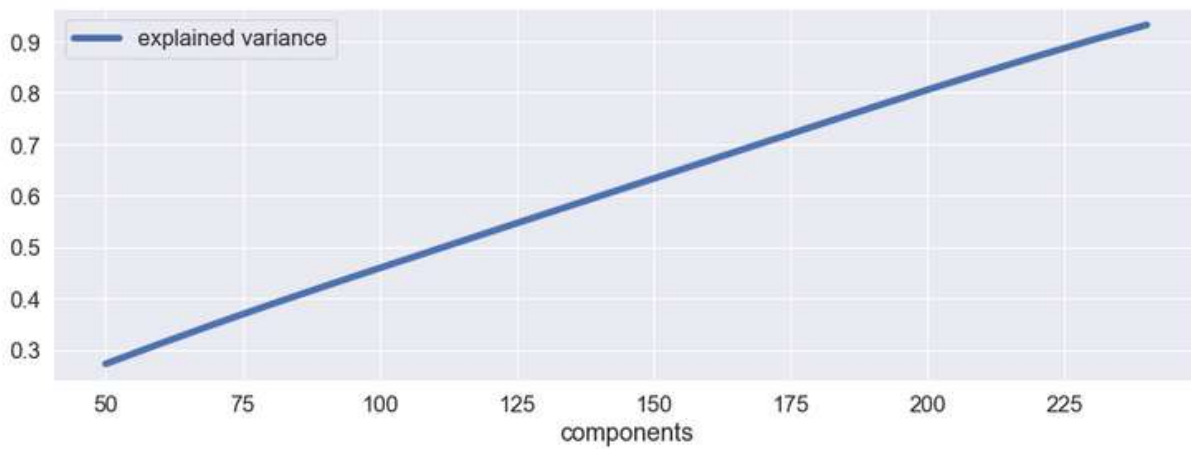
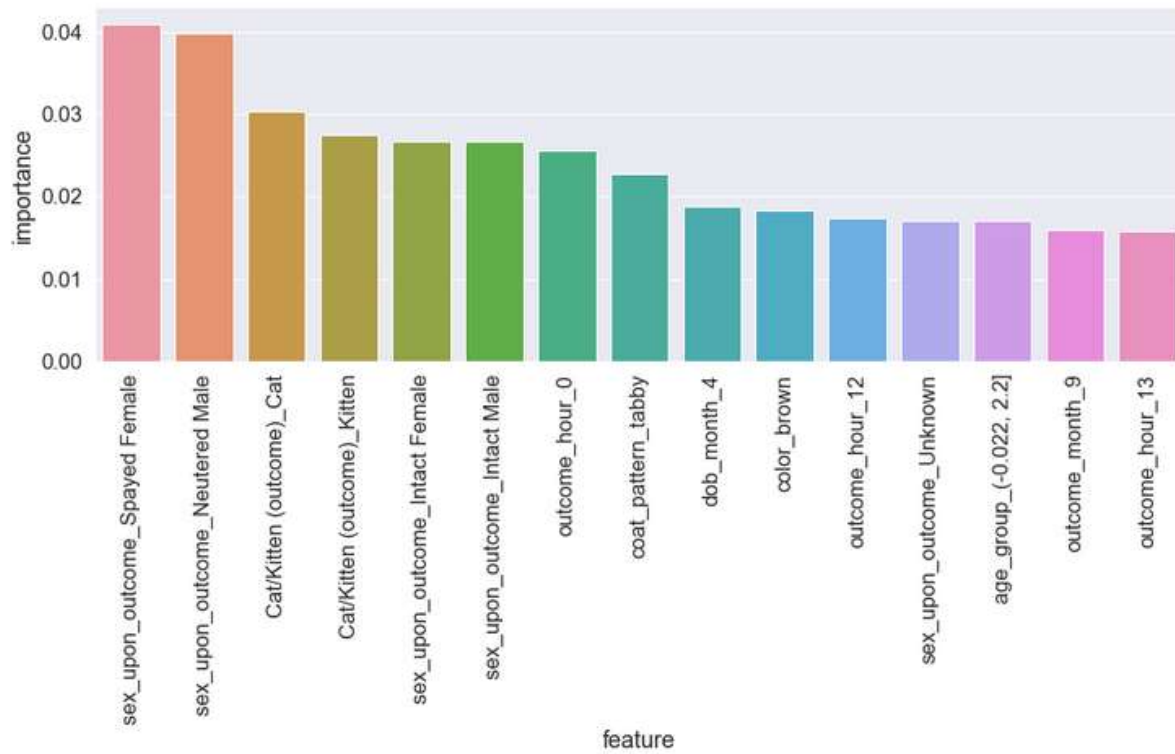


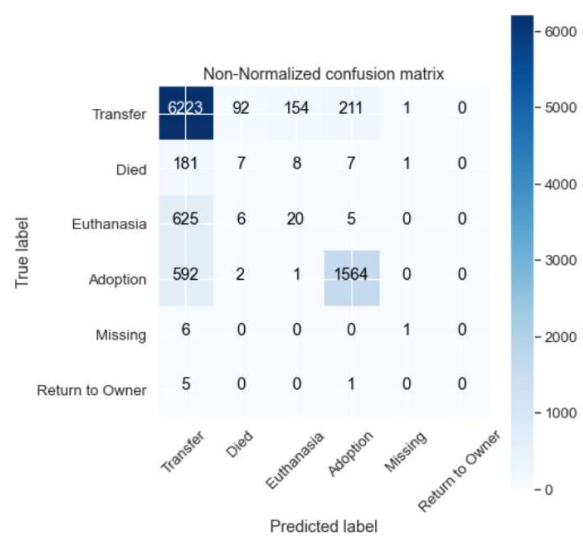
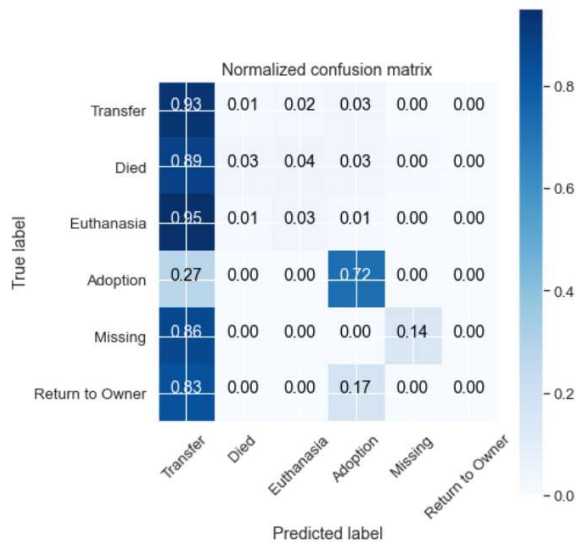




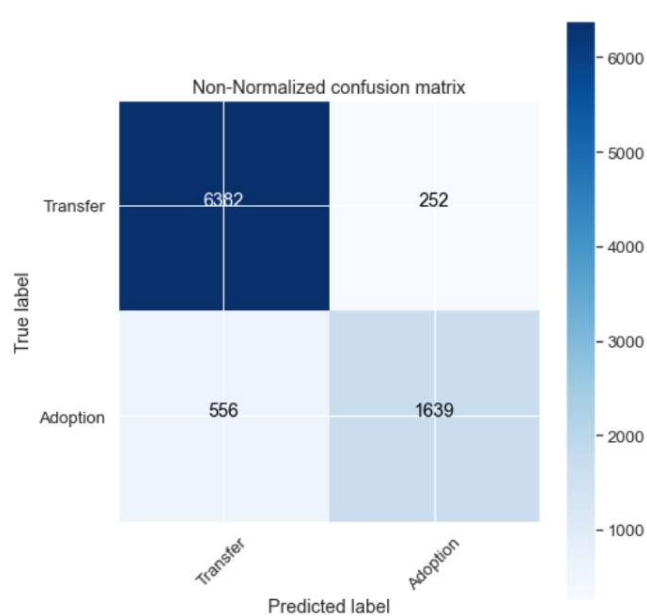
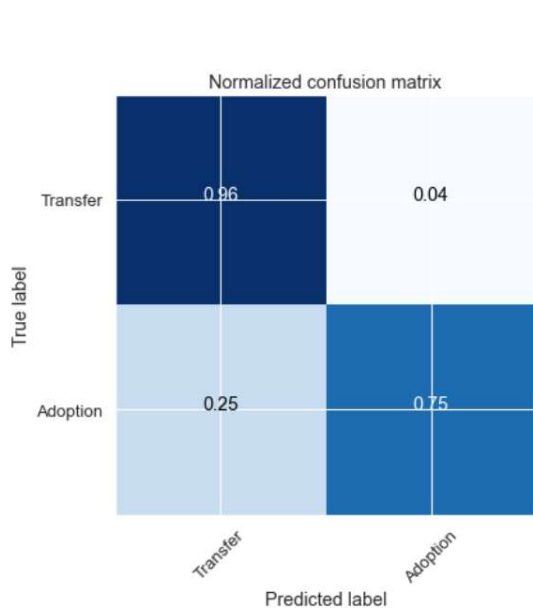
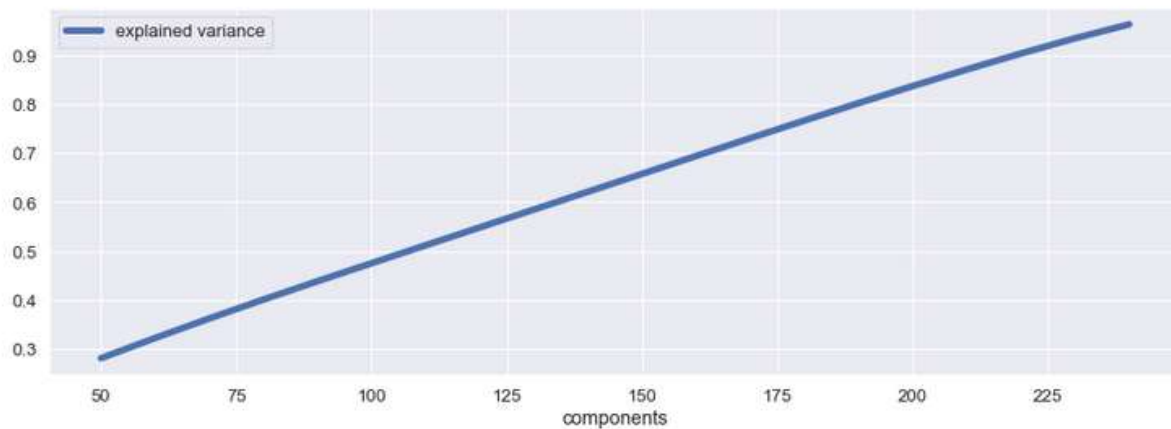
## Machine Learning Model Visualizations & Outputs

### Cats:



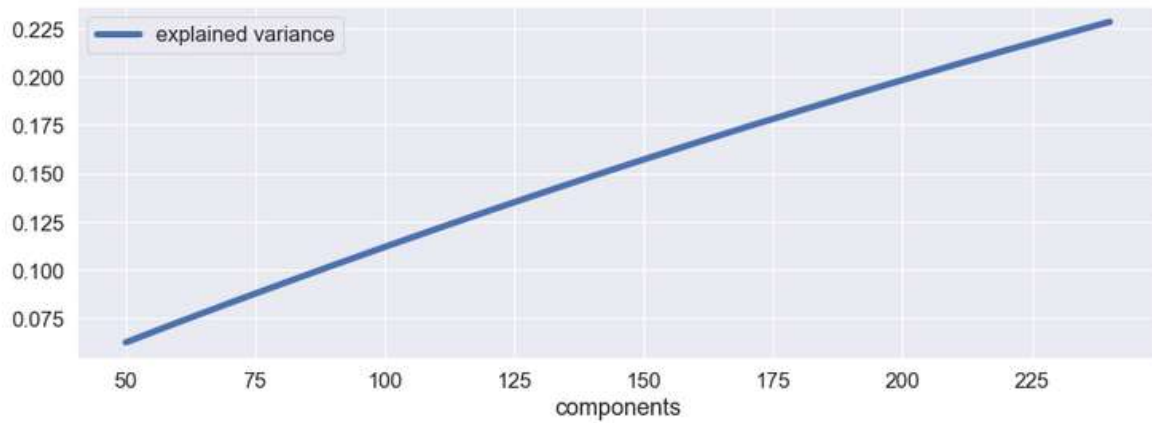
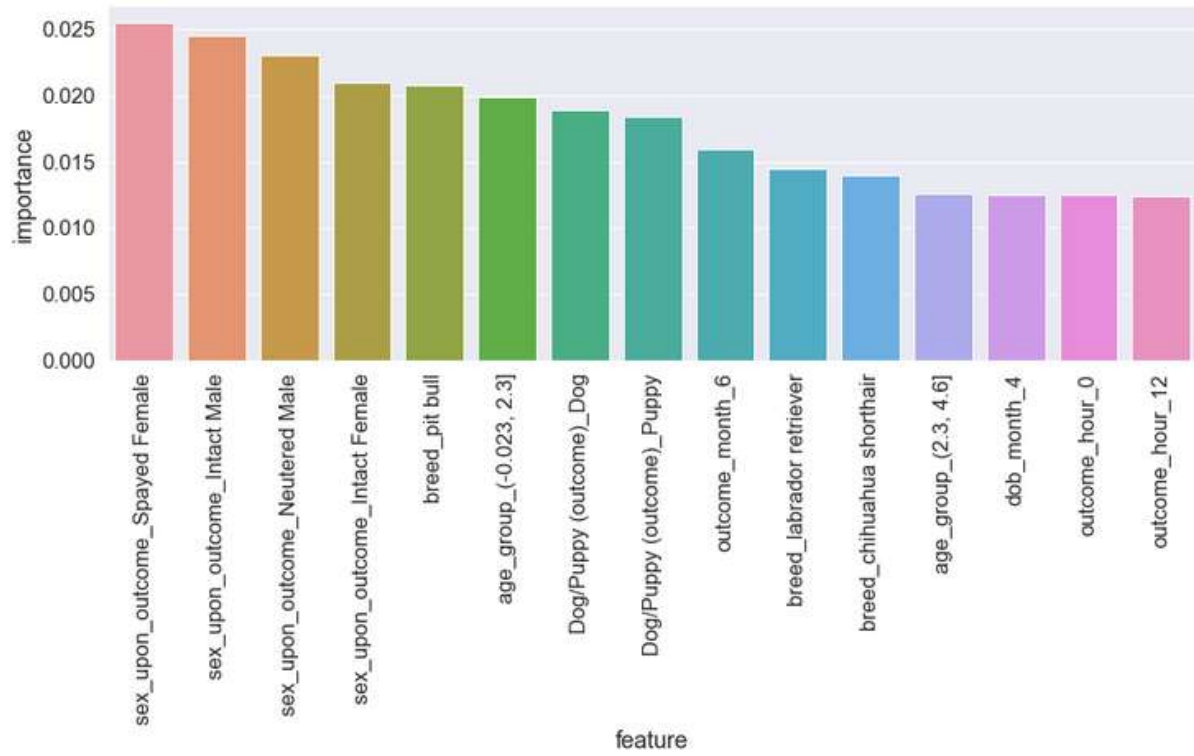


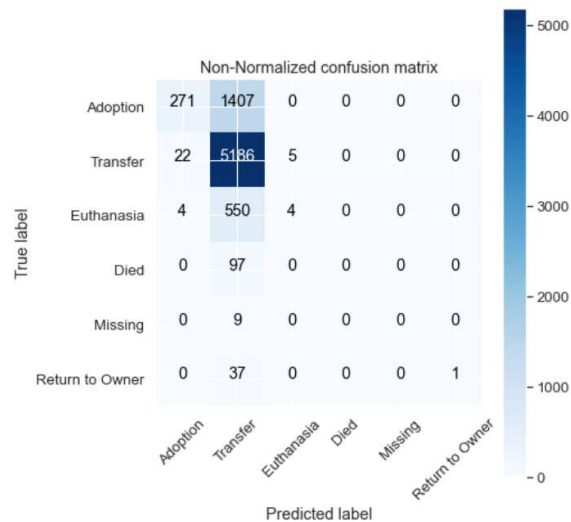
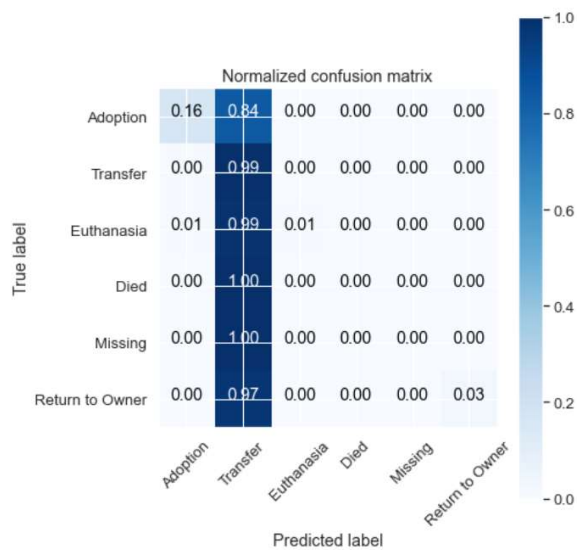
	precision	recall	f1-score	support
Transfer	0.82	0.93	0.87	6681
Died	0.07	0.03	0.05	204
Euthanasia	0.11	0.03	0.05	656
Adoption	0.87	0.72	0.79	2159
Missing	0.33	0.14	0.20	7
Return to Owner	0.00	0.00	0.00	6
accuracy			0.80	9713
macro avg	0.37	0.31	0.33	9713
weighted avg	0.76	0.80	0.78	9713



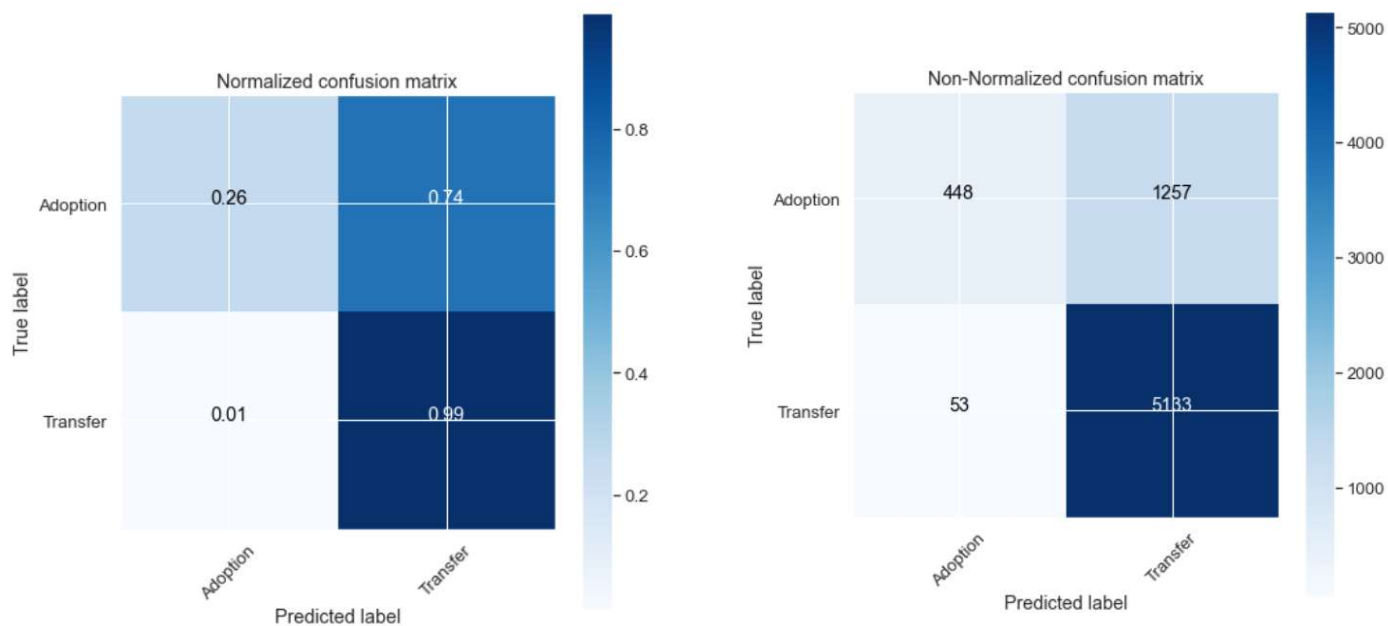
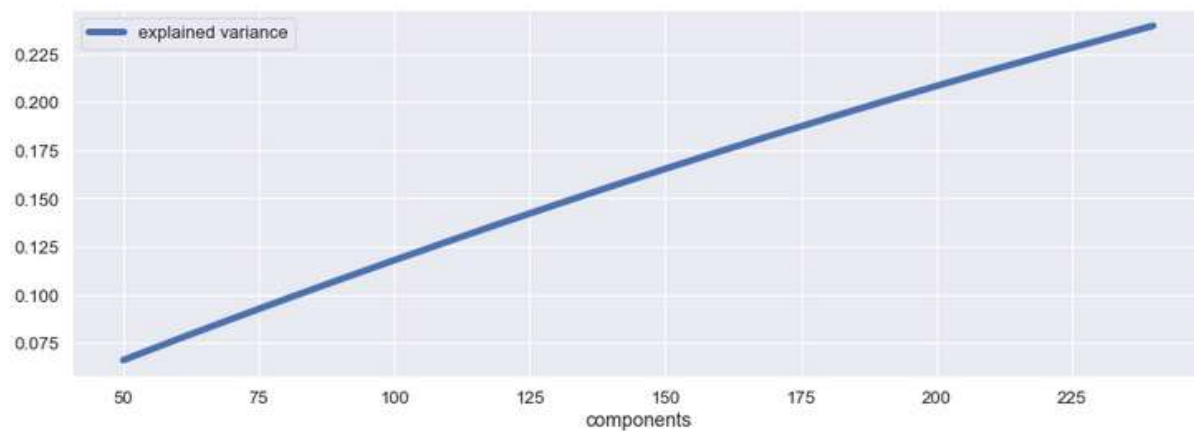
	precision	recall	f1-score	support
Transfer	0.92	0.96	0.94	6634
Adoption	0.87	0.75	0.80	2195
accuracy			0.91	8829
macro avg	0.89	0.85	0.87	8829
weighted avg	0.91	0.91	0.91	8829

### Dogs:





	precision	recall	f1-score	support
Adoption	0.91	0.16	0.27	1678
Transfer	0.71	0.99	0.83	5213
Euthanasia	0.44	0.01	0.01	558
Died	0.00	0.00	0.00	97
Missing	0.00	0.00	0.00	9
Return to Owner	1.00	0.03	0.05	38
accuracy			0.72	7593
macro avg	0.51	0.20	0.19	7593
weighted avg	0.73	0.72	0.63	7593



	precision	recall	f1-score	support
Adoption	0.89	0.26	0.41	1705
Transfer	0.80	0.99	0.89	5186
accuracy			0.81	6891
macro avg	0.85	0.63	0.65	6891
weighted avg	0.83	0.81	0.77	6891