# Group Project

# Animal Adoptability Prediction

**By**

**Dhaval Shah (802973719)**

**Pritesh Pimpale (803000983)**

**Dipika Mahashabde (803004274)**

**CPSC 483**

## Professor: Kenytt Avery

**Department of Computer Science**

**California State University, Fullerton**

**Spring 2016**

# CONTENTS

# 1   INTRODUCTION

Kaggle Competition it is a platform where different users take part, analyze the data and perform different algorithms, the prediction produced by the algorithm is submitted on the online site and the scores are produced by the competition. The results have ranks and scores, the scores are important of all if you have a score near to 0 it's a really good prediction produced by your algorithm. They have rules and regulation before submitting the predicted outputs.

We had to choose a competition which was active on Kaggle. Therefore we looked at different competitions and found one which was on Shelter Animals (https://www.kaggle.com/c/shelter-animal-outcomes).

Shelter Animal outcomes, the data for this competition is taken from Austin Animal Center. Every year many animals are given up by owners or lost or even go through cruelty situation. And leaving this there are many animals which get euthanized. Therefore the Austin animal center wanted to predict the outcome if the animals will go through transfer to other shelter or return to owner or die of natural causes or get euthanized or get adopted.

In simple words we have classify if the animal will get

1. Transferred
2. Return to owner
3. Died
4. Adopted
5. Euthanized

Data in data files

1. **AnimalType**          : Types of animals (Cat or Dog)
2. **Name**                : Name of the animal(pets)
3. **DateTime**            : Date and time at which the outcome was measured
4. **OutcomeType**         : The result which we have to predict on test data
   a. Transferred
   b. Return to owner
   c. Died
   d. Adopted
   e. Euthanized
5. **OutcomeSubType**      : the reason for the outcome mentioned in the above column
6. **SexuponOutcome**      : Gender of the animal
   a. Intact Female
   b. Spayed Female
   c. Intact Male
   d. Neutered Male
   e. Unknown
7. **AgeuponOutcome**      : age of the animal(in days/week/months/years)
8. **Breed**               : Breed of particular animals

**9. Color** : Color of the animals(may one or two colors are given)

**Data:**

# 2 PREPROCESSING OF DATA

Preprocessing of data is necessary to run the data mining algorithms on the given data. As many of the columns have data which gives the multiple values of the attribute. e.g. Color column has more than one color for the mixed colored cats and dogs. So this type of data is necessary to be preprocessed before attempting any algorithm on it.

- **Name**
  Animals have some name given to them. This name can be a feature for the adoption of the pet by new owner. There are many animals who don't have names so we have kept it blank.
- **Animal Type**
  This attribute has value either 'cat' or 'dog'.
- **sex upon outcome**
  This attribute stores the sex of the pet animal. This attribute has values
  a. Intact Female
  b. Spayed Female
  c. Intact Male
  d. Neutered Male
  e. Unknown
- **Name Frequency**
  We have also calculated the frequency of all the names which were registered
- **Has Name**
  We also created a new data which stores the value 'Yes' if the name is present and 'No' if the name is not present. This parameter is much more significant in case of cats. For cats the adoption rate is much higher if a cat has name.
- **Date Time**
  Date Time attribute stores the time of the outcome of the train dataset. This attribute we have divided into many columns which are
  1. Year
  2. Month
  3. Day
  4. Hour
  5. Minutes
- **Age upon Outcome**
  The Age was in all different forms like days weeks months and years, which would make it difficult to analyze the data. So we convert the in format which is we convert the data into weeks.
- **Age Groups**
  Age is one of the factor for the decision of the owner to adopt a pet. So instead of keeping the data as it is we converted age into different age ranges. So we were able to group the age into ranges which provided an overview of the data.
- **Breed**
  We have different types of breed for each animals so we have checked if the animals (dog/cat) are of pure breed or mixture of breeds.

3

- **Mix Breed**

  The breed attribute has the value Mix if the if the breed is mixed. We extracted this value into new column to set if the breed is mix or not. If the animal is of pure breed we assign it 'normal' and if it is mixture we would assign 'mix'.
- **Color**

  We have different colors for all the animals so to see if it has mixture of two colors or a single color. All values of the color are stored in the color column.
- **Mix color**

  We have added one more column to see if it has two colors or no. If multiple colors are there for animal we set the Mix color attribute to true and false if it is single color.

After processing the data we have converted the data into numeric form to use it in the algorithms.

e.g.

For Animal Type
- Dog = 0
- Cat = 1

# 3 DATA ANALYSIS

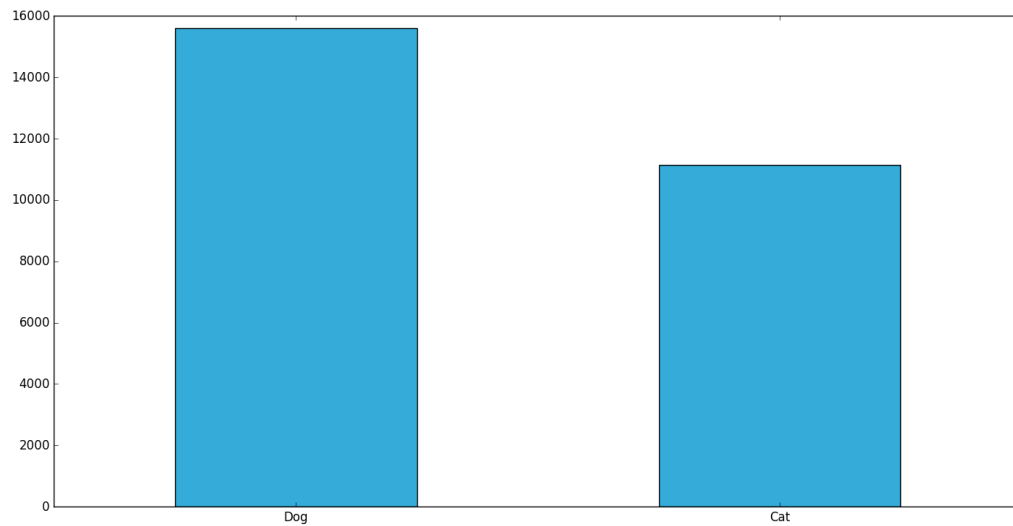## 3.1 ANIMAL TYPE DISTRIBUTION



Fig. Cat and Dog

After seeing the graph we get to know that in the train data we have more dogs than cats.
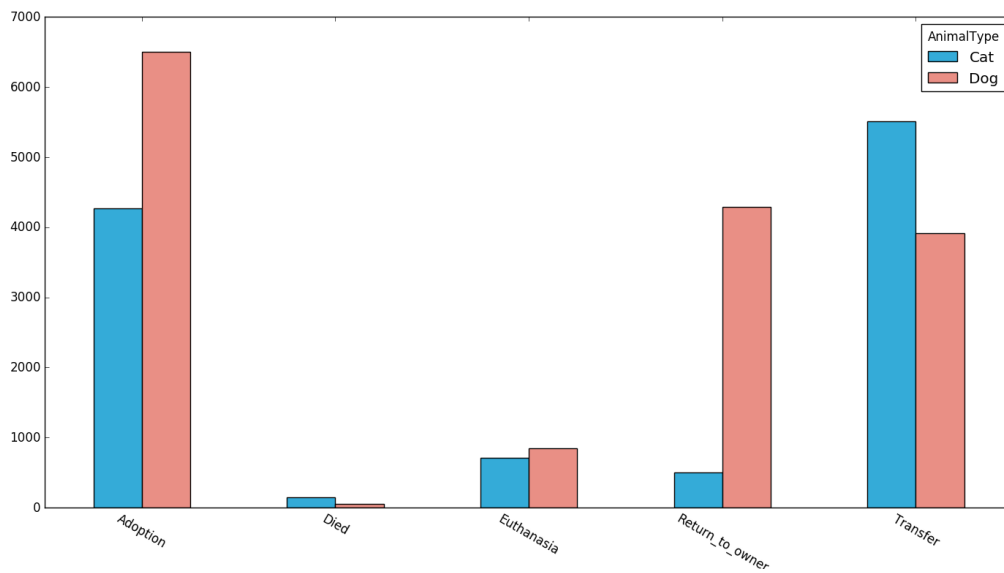
## 3.2 OUTCOME DISTRIBUTION



Fig. Animal type vs outcome

This shows that Adoption of animals is more than other outcome. Also dog has more probability of getting adopted or going return to owner than the cats.
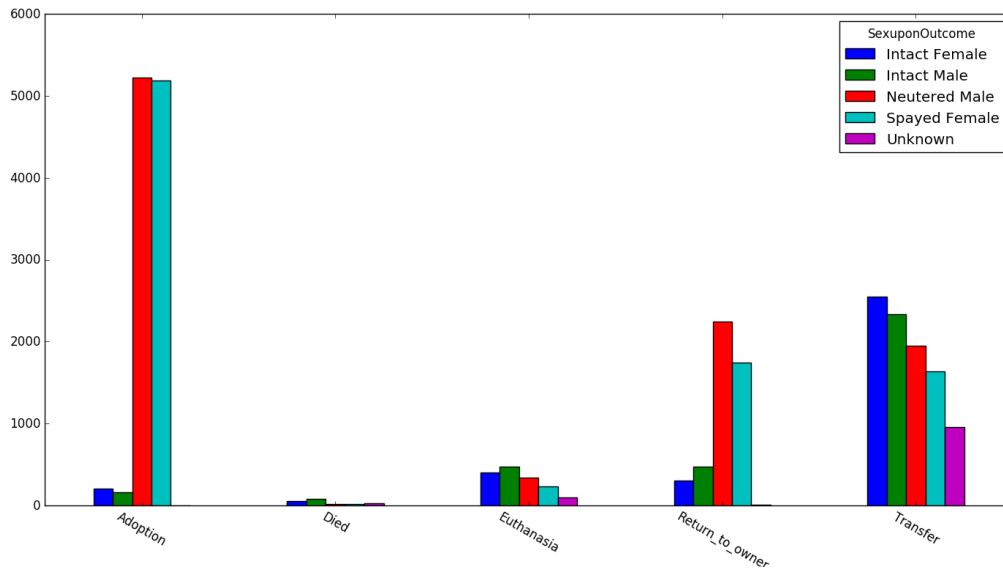
## 3.3  SEX VS ANIMAL TYPE



Fig. sex vs animal type

This data shows that adoption of neutered male and spayed female is more likely to happen than other types of the sex. Also same trend is seen for the return to owner outcome. For the transfers the sex doesn't seem to matter much.

- **Adoption:** Neutered Male and Spayed Female are most of time adopted.
- **Died:** sex doesn't matter much for the outcome of the pet.
- **Euthanasia:** intact female and intact male are more euthanized more than neutered or spayed. But still this difference is not much and doesn't provide any significance.
- **Return_to_owner:** Neutered male and spayed female are returned more than others.
- **Transfer:** Intact male and female is more than others. But data doesn't show much significant difference between them.
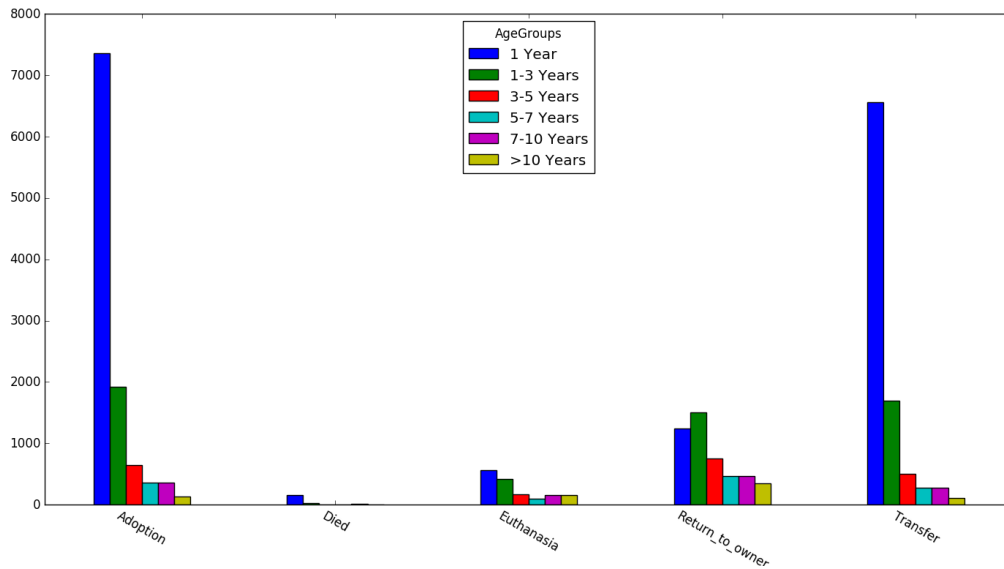
## 3.4  AGE GROUPS VS OUTCOME TYPE



Fig. age group vs outcome Type

Adoption and transfer of the animals is more for the animals aged less than 1 year than for the animals who are old.
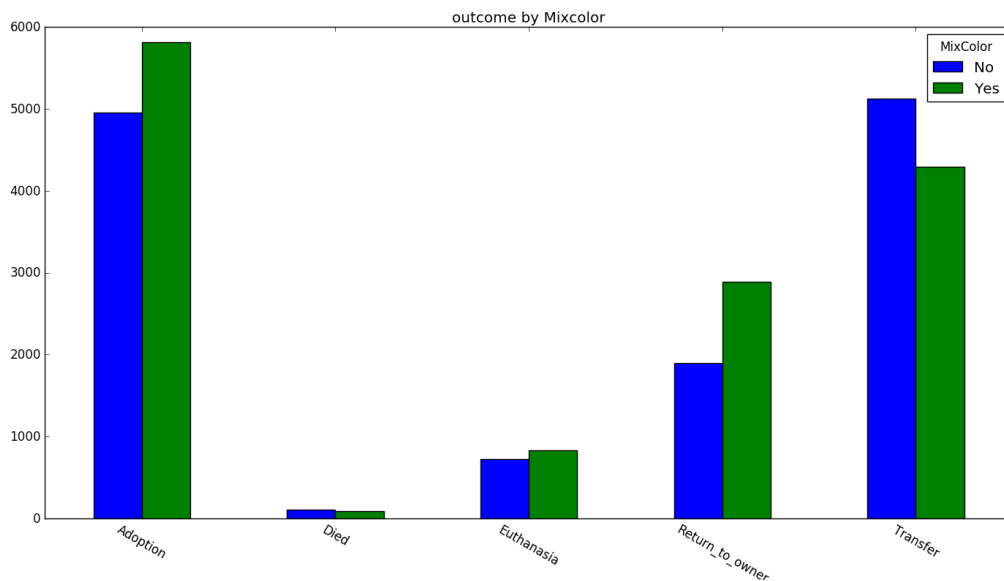
## 3.5  MIX COLOR VS OUTCOME TYPE



Fig. Mix Color vs Outcome Type

This graph tells us that how many animals have mix color and how many have only one color. This graph doesn't give any significant information as for all the outcome data is pretty similar.
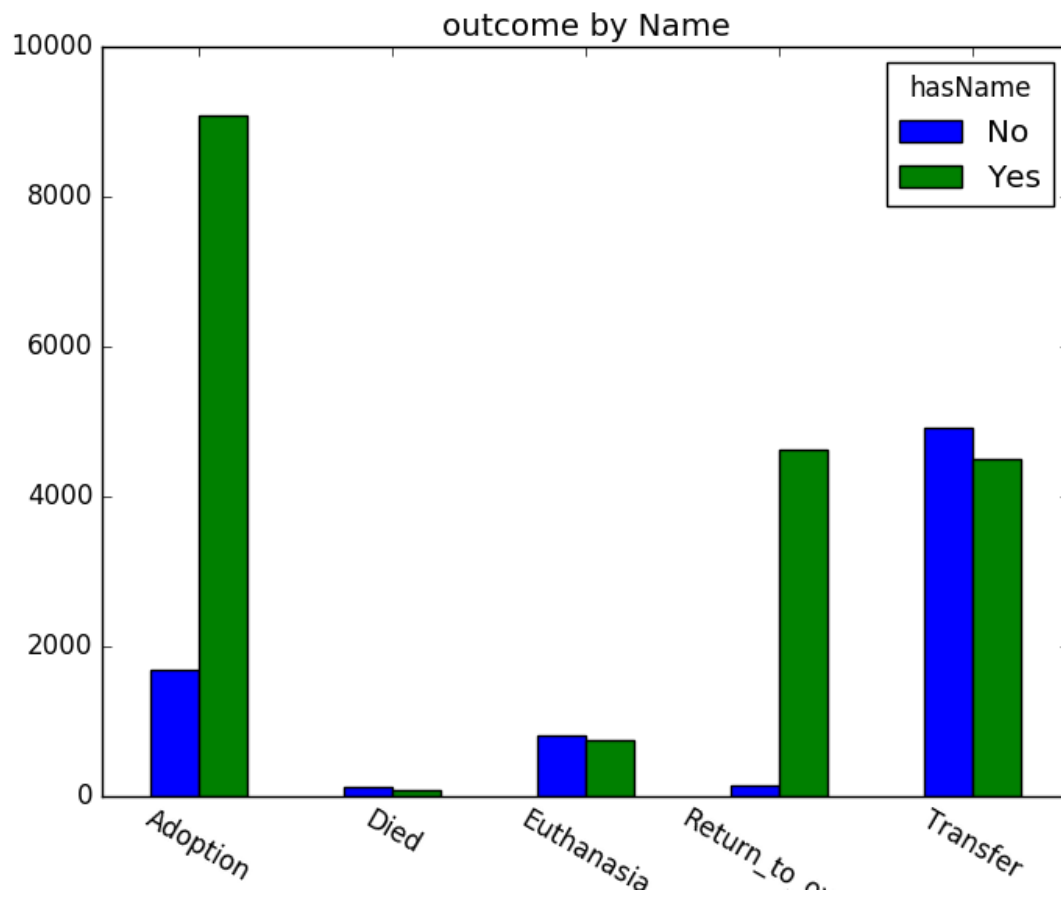
## 3.6 HAS NAME VS OUTCOME TYPE



Fig. has Name vs outcome type

This graph shows that having a name to a pet animal is significant factor for animal getting adopted or returning to owner. While having name doesn't matter much for rest of the outcomes.
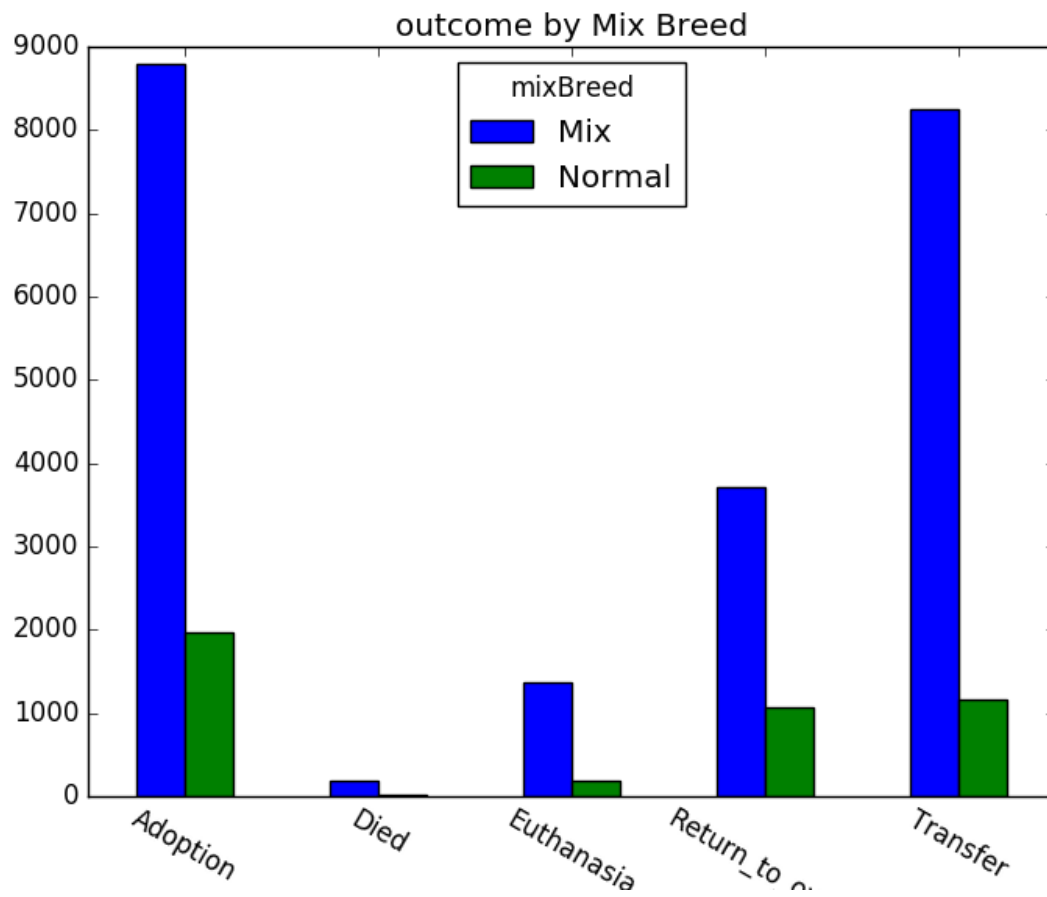
## 3.7 MIX BREED VS OUTCOME TYPE



Fig. Mix Breed vs Outcome Type

Mix breed is seen more in all the outcomes. This shows that the pet animals with mix color are more in the dataset, so this also doesn't provide much information for predicting the outcome.

# 4   RESULTS OBTAINED FROM TEST AND TRAINING DATA

- **Naïve Bayes**

Naive Bayes is a simple technique for constructing classifiers: models that assign class labels to problem instances, represented as vectors of feature values, where the class labels are drawn from some finite set. It is not a single algorithm for training such classifiers, but a family of algorithms based on a common principle: all naive Bayes classifiers assume that the value of a particular feature is independent of the value of any other feature, given the class variable. For example, a fruit may be considered to be an apple if it is red, round, and about 10 cm in diameter. A naive Bayes classifier considers each of these features to contribute independently to the probability that this fruit is an apple, regardless of any possible correlations between the color, roundness and diameter features.

We have implemented naïve Bayes algorithm. We used Naive Bayes from SK-learn (scikit learn, n.d.) Library.

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of independence between every pair of features.

Gaussian NB implements the Gaussian Naive Bayes algorithm for classification. The likelihood of the features is assumed to be Gaussian:
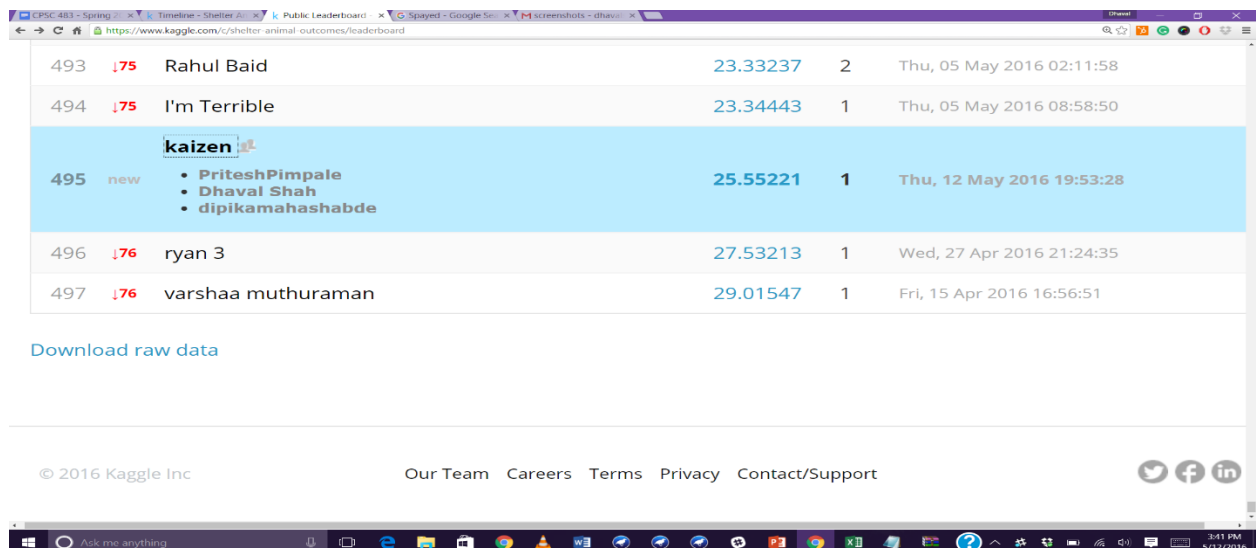
$$P(x_i \mid y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} \exp\left(-\frac{(x_i - \mu_y)^2}{2\sigma_y^2}\right)$$

The parameters $\sigma_y$ and $\mu_y$ are estimated using maximum likelihood.

We have train the data on the outcome type and then predicted whether the outcome type is Adoption, Transferred, Return to owner, Died, Adopted, and Euthanized.

Our initial approach was to use this classifier for the prediction. The initial results were not that good for the prediction as there were no extra parameters considered for the prediction.

Kaggle score is as follow.

- **Random Forest**
  Many Classification trees are grown by using random forest (scikit-learn, n.d.). An input vector is given to predict its class. This input vector passes through all tress which are produced. Each tree gives a result that is, which class it belongs to. Votes are considered from all the trees to decide the final outcome of the classifier. The class with most vote is assign to the input given.
  **Growing of the trees**
  There are 'N' number of observation in the training set. Every tree selects observations at random by replacing the actual training data. These new data generated are used for the growing of the tree. There are 'M' number of attributes, they are also selected at random, and the best split on these attributes (selected) is used to split the node. The value of the M can be given while running the algorithm. While growing the tree the selected attribute won't change till the tree is grown to the output leaves. (random forest, n.d.)

  Our next approach was to use random forest classifier for the prediction.
  We tried with the original parameters and following configurations for the random forest classifier.
  - No. of trees = 10
  - max features = 5
  - max depth = 10

  Initial result with this parameters is as below



Random forest algorithm can give better predictions if number of trees are set to a larger values.

Also the new derived parameters are used while training the random forest classifier.

Parameters used for the final submissions are

1. Name
2. NameFreq
3. hasName
4. AnimalType
5. SexuponOutcome
6. AgeuponOutcome
7. AgeGroups
8. Breed
9. mixBreed
10. Color
11. MixColor
12. firstColor
13. Year
14. Month
15. Day
16. Hour
17. Minute

For random forest classifier there is no need to do the feature selection as it is best trained with all the features provided. Each tree generated uses only 'n' number of features specified while running it from all features. It also trains the generated trees on the result of the earlier trees. This increases the prediction accuracy of the trees in the classifier.

We tried with multiple values for the number of trees and number of features. The final submission was the one with minimum logloss obtained from the training data.

Final parameters used for the prediction.

- o No. of trees = 1000
- o max features = 10
- o max depth = 10

**Final result from Kaggle**

# 5 ENVIRONMENT

**Environment**: - Python (install the new version)

**Libraries: -**

1. Pandas: - It offers data structures and operations for manipulating numerical tables and time series
2. Scipy: - It helps in scientific and technical computing.
3. Sklearn: - They have all algorithm with their functions.
4. Numpy: - It helps in mathematical computation.
5. Matplotlib: - This library is used to plot graphs in python, it has inbuilt functions for it.

 **Tool: -** Visual Studio Code.

1. **Installation of Python and it libraries**
   **Step 1.**
   Download python: - **https://www.python.org/downloads/**
   **Step 2.**
   Install all the libraries needed
   1. Install pip(This will help us install other libraries)
   2. Install scipy (pip install scipy)
   3. Install numpy (pip install numpy)
   4. Install sklearn (pip install sklearn)
   5. Install pandas (pip install pandas)
   6. Install matplotlib (pip install matplotlib)


2. **Download the data from Kaggle**
   Data: - https://www.kaggle.com/c/shelter-animal-outcomes/data (Kaggle , n.d.)
   - Download all the files on this link.
   - Train.csv: - this data will help us to model the data using the algorithm.
   - Test.csv: - the modeled algorithm will help us to predict on the test data.
   - Sample_Submission.csv: - the predicted values need to be stored in the format shown in this file.

# 6 REFERENCES

*Kaggle* . (n.d.). Retrieved from Kaggle : https://www.kaggle.com/c/shelter-animal-outcomes/data

*random forest*. (n.d.). Retrieved from random forest:
https://www.stat.berkeley.edu/~breiman/RandomForests/cc_home.htm#intro

*scikit learn*. (n.d.). Retrieved from scikit learn: http://scikit-learn.org/stable/modules/naive_bayes.html

*scikit-learn*. (n.d.). Retrieved from scikit-learn: http://scikit-learn.org/stable/modules/ensemble.html#forest