

# Investigating Factors into Squirrels' Attitudes towards Humans in New York

Team BCHZ - Nick Bayer, Richard Cui, Laura Han, Anna Zhang

2023-12-01

## Introduction and Data

As a result of the continuous human development, animals are inevitably interacting with humans more often. However, this form of interaction has mostly shown to be a disturbance to animals [1]. Animals see humans as a threat, so it is no surprise that they would treat the presence of humans the same way they would when they face other predators. Nevertheless, recent studies show that the squirrels actually act differently, as characterized by a phenomenon called synurbization, or the process of becoming urbanized [2].

In an effort to investigate these two competing theories, and to better understand the dynamic between squirrels and humans, we carry out this research project to explore what factors affect whether a squirrel is indifferent to human presence. From there, we would like to deduce whether the squirrels' attitude to humans are caused by human presence or other factors such as their species.

We hypothesize that the age category, location, distance above ground when spotted, number of activities that the squirrel was doing, sound that the squirrel makes, and whether squirrel is disturbed by human activities (as measured by features like approaching or running away from humans and tail signs) could have a relationship with the attitude of the squirrel (whether indifferent or not).

## Data description

We are sourcing our data set from the TidyTuesday project on GitHub. Their data originally came from The 2018 Squirrel Census, a project based on the sightings of squirrels in Central Park, New York City.

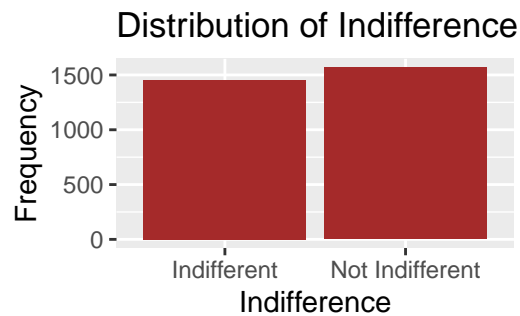
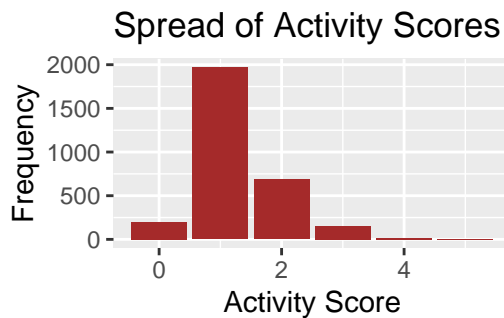
In October of 2018, the Squirrel Census Team and a group of over 300 volunteers collected the data based on squirrel sightings around Central Park.

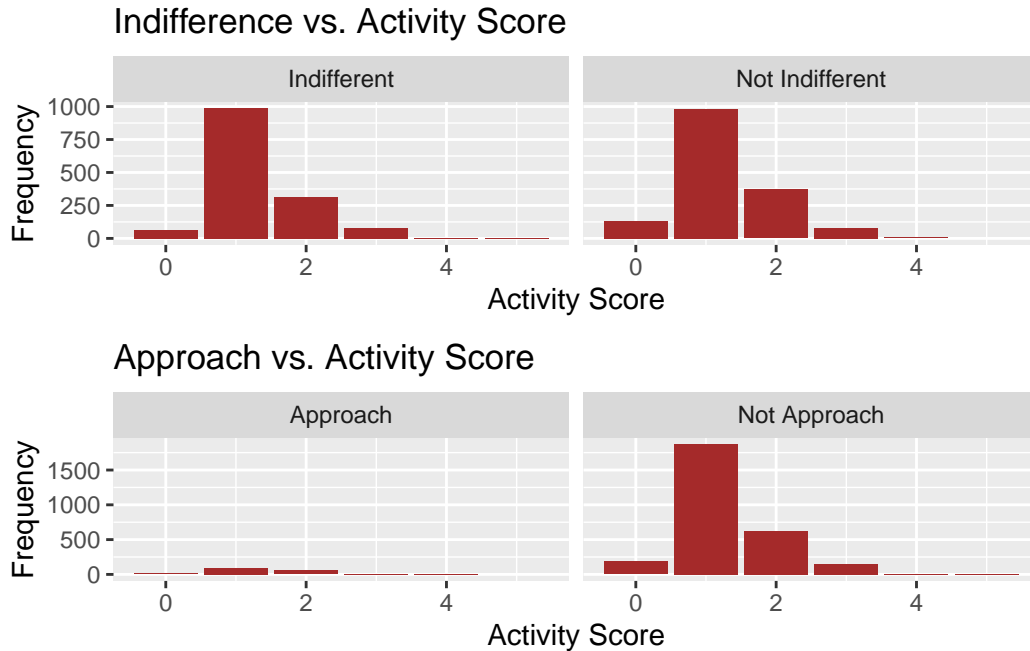
The dataset has 3023 observations and 31 variables, and it gives a wide range of observations and characteristics. It first gives us the location, in both longitude and latitude, the hectare of the park the squirrel was located in, the date, and whether it was found in the AM or PM. It also assigns each squirrel a unique ID. It has information on whether the squirrel is an adult or a juvenile, its primary and highlight fur colors, and has a number for the sequence of sightings in one session. It also contains data on their exact location, their distance from the ground, and the objects they were found on. There is data for the activities the squirrel was found doing, ranging from running to foraging, with a separate column for any activity that was not chosen to be a column. It gives data on the sounds the squirrel made and tail movements, if any. Finally, it has 4 columns for the squirrel's response when facing humans, being either that the squirrel approached, was indifferent, ran away, or any other action.

## Initial exploratory data analysis

We create a numeric variable named `Activity_Score` that encapsulates the number of activities a squirrel is engaged in during the span of observation. The distribution of the `Activity_Score` variable is unimodal, slightly right skewed, and has a median of 1. This means that most squirrels in the data set were only engaged in 1 activity during the span of observation.

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
0.000	1.000	1.000	1.278	2.000	5.000





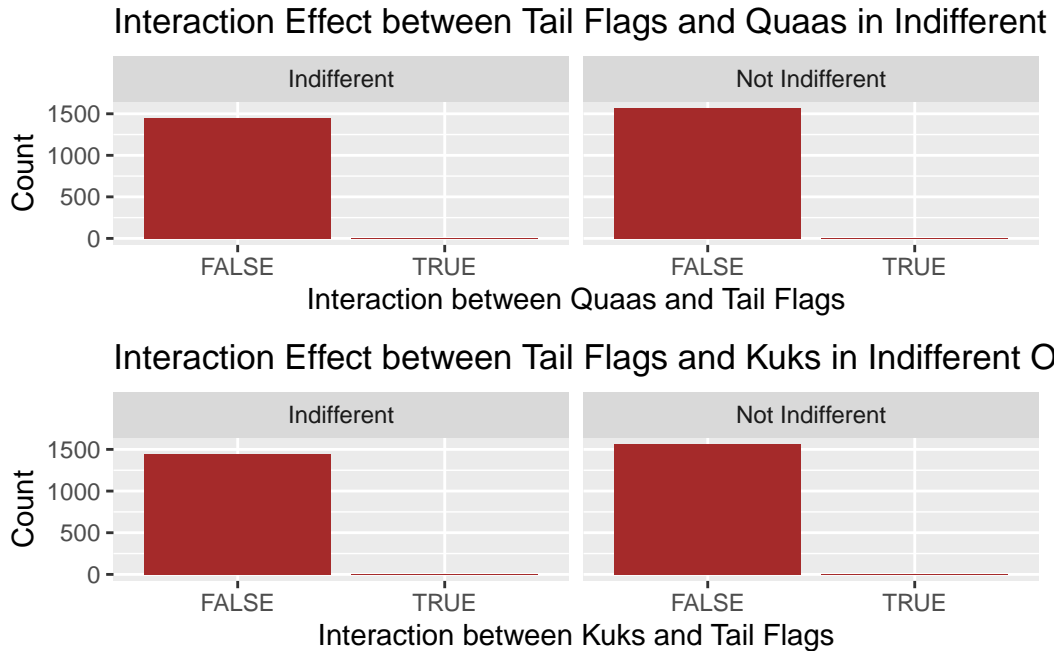
There seems to be roughly equal number of squirrels that were and weren't indifferent to the humans.

There were far fewer squirrels that approached humans than squirrels that didn't. This could show that they are not as synurbanized as we thought.

The distribution of `Activity_Score` for squirrels that were indifferent vs. not indifferent is roughly the same shape. Therefore, activity score and indifference do not seem to be correlated.

Since there were so few squirrels that approached the researchers, it is hard to tell whether the graphs are very different in shape. The difference may be something that we explore further.

An interaction that we would like to examine further is whether the presence of tail activities, specifically `tail flags`, when interacting with `Quaas/Kuks`, influences the likelihood that a squirrel is indifferent to humans. When squirrels flag their tails, the motion exaggerates their size and is used to confuse rivals or predators. When squirrels are heard quaaing, they are making an elongated vocal communication to indicate the presence of a ground predator. Meanwhile, kuk is a chirpy vocal communication used for a variety of reasons. Therefore, seeing the interaction between tail flags and quaas, more broadly, kuks, could help us understand whether or not the squirrels are indifferent to humans or not. In both graphs below, there are more counts of false interaction between both `Tail Flags` and `Quaas` as well as `Tail Flags` and `Kuks` (1600) than true interaction (1400), which shows a mild relationship between the variables.



## Data dictionary

The data dictionary can be found [here](#).

## Methodology

### Analysis approach

Our response variable is **Indifferent**, which is a categorical variable that indicates whether or not the squirrel is indifferent to humans. Potential predictors include **Activity\_Score** (a quantitative variable that records the number of activities the squirrel is observed doing), **Age** (a categorical variable that indicates whether the squirrel is adult or juvenile), **Fur\_Color** (categorical variable), **Location** (categorical variable), **Above\_Ground\_Measurement** (quantitative variable), sounds that the squirrels are making (categorical variables including **Kuks**, **Quaas**, and **Moans**), **Tail\_flags** and **Tail\_twitches**, which are also categorical.

To explore the relationship between whether or not the squirrel is indifferent and the predictor variables, such as age category, location, distance above ground when spotted, number of activities the squirrel is observed doing, sound that the squirrel makes, and their tail signs, we plan to use logistic regression. We are using logistic regression because our dependent variable, **Indifferent**, is categorical. We will compare logistic regression models using AIC and BIC

to evaluate what predictor variables and what interactions between the variables should be included in the model to best predict the attitude of the squirrel towards humans. We will also perform 10-fold cross-validation for model comparison.

## Model 1

Model 1 includes all variables stated in our hypothesis.

Recipe for Model 1 steps: 1. Change response variable into factors. 2. Map all “FALSE” in **Above Ground Sighter Measurement** variable to 0, then convert the variable type to integer. 3. Create dummy variables for all nominal predictors. 4. Create interaction terms between **Quaas** and **Tail flags** 5. Remove all variables with zero variance.

.metric	.estimator	mean	n	std_err	.config
accuracy	binary	0.533	10	0.010	Preprocessor1_Model1
roc_auc	binary	0.537	10	0.011	Preprocessor1_Model1

mean_aic	mean_bic
2765.376	2843.233

## Model 2

As the EDA has shown, **Activity\_Score** does not seem to have a relationship with **Indifferent**, so we take it out for Model 2. We conduct a drop-in-deviance test that confirms that the coefficient of **Activity\_Score** is not statistically significant from 0 because the p-value of the test  $0.72 > 0.05$ .

Recipe for Model 2 steps: 1. Change response variable into factors 2. Map all “FALSE” in **Above Ground Sighter Measurement** variable to 0, and not “FALSE” to 1, then convert the variable type to factors. 3. Create dummy variables for all nominal predictors 4. Remove all variables with zero variance.

.metric	.estimator	mean	n	std_err	.config
accuracy	binary	0.524	10	0.014	Preprocessor1_Model1
roc_auc	binary	0.539	10	0.012	Preprocessor1_Model1

mean_aic	mean_bic
2762.821	2829.475

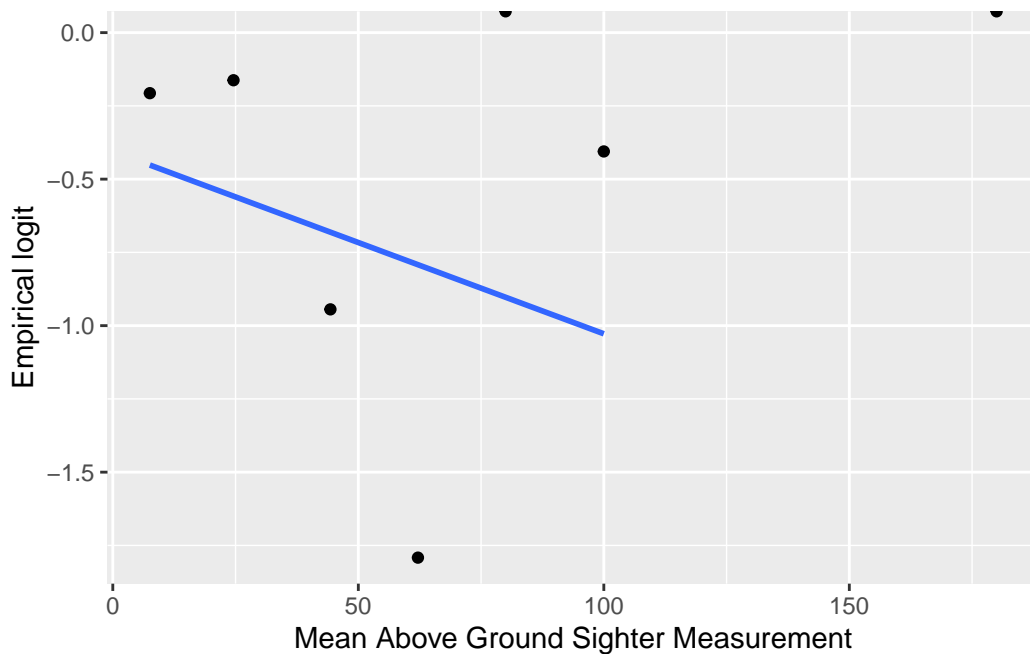
Drop-in-deviance p-value
0.7227502

## Model Conditions

We will check model conditions for Model 2 since this is the better performing model due to higher accuracy and AUC and lower AIC and BIC.

### Linearity

There is one numeric variable in Model 2, Above Ground Sighter Measurement.



The linearity condition is not satisfied. There is not a linear relationship between the empirical logit and the predictor variable of Above Ground Sighter Measurement. We could fix this condition if we had more specific data with fewer NA values. Additionally, since we converted the false values to zeros, that may have led to the variable not having a linear relationship with the empirical logit. Because the condition is not satisfied, the coefficients and p-value for Above Ground Sighter Measurement in our model may not be accurate, and we are unable to conclude if Above Ground Sighter Measurement is a significant predictor.

## Randomness

The data was collected from the sightings of squirrels in Central Park, NYC from a group of volunteers. Although the volunteers are not randomly sampled, the sample of squirrels can be considered as random since we do not have reason to believe that the squirrels collected in this study differ systematically from squirrels in the rest of the world in regards to their indifference to humans. Therefore, randomness is satisfied.

## Independence

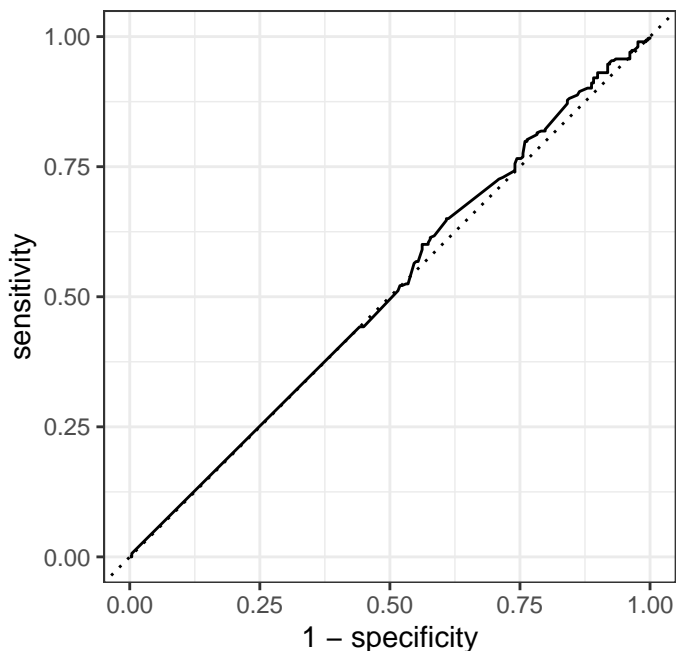
The data are not spatially or time correlated since the data were collected in one time snapshot and in one location. Therefore, we have no reason to believe that independence was violated.

Thus, all model conditions for logistic regression is satisfied.

## Results

Since Model 2 has a higher accuracy and AUC, as well as lower AIC and BIC, we will choose Model 2. Next, we fit the model to the entire squirrels dataset and interpret the coefficients.

term	estimate	std.error	statistic	p.value
(Intercept)	-0.108	1.438	-0.075	0.940
Above Ground Sighter Measurement	-0.011	0.005	-2.357	0.018
KuksTRUE	-0.542	0.262	-2.065	0.039
QuaasTRUE	-0.908	0.444	-2.045	0.041
MoansTRUE	-13.221	535.411	-0.025	0.980
Tail flagsTRUE	-0.044	0.196	-0.223	0.823
Tail twitchesTRUE	-0.018	0.120	-0.149	0.881
Age_Adult	-0.102	1.421	-0.072	0.943
Age_Juvenile	-0.414	1.426	-0.290	0.772
Primary Fur Color_Cinnamon	0.133	0.251	0.530	0.596
Primary Fur Color_Gray	0.272	0.225	1.208	0.227
QuaasTRUE_x_\'Tail flags\'TRUE\'	14.645	262.751	0.056	0.956



```
# A tibble: 1 x 3
  .metric .estimator .estimate
  <chr>   <chr>       <dbl>
1 roc_auc binary      0.511
```

	Truth	
Prediction	FALSE	TRUE
FALSE	113	124
TRUE	145	179

Of all the variables we examined, `KuksTRUE`, `QuaasTRUE`, `Above Ground Sighter Measurement` were the only terms that had coefficients with significant p-values. The odds that a squirrel is indifferent to a human is multiplied by a factor of 0.582 ( $\exp(-0.542)$ ) if it kuks and 0.403 ( $\exp(-0.908)$ ) if it quaas, holding all else constant. For every additional 1 meter in the squirrel's location above ground, the odds that it is indifferent is multiplied by a factor of 0.989 ( $\exp(-0.011)$ ), holding all else constant. In other words, kuks, quaas, and being above ground decrease the odds that a squirrel is indifferent to humans. However, the interaction term between `KuksTRUE` and `Tail flagTRUE` does not have a significant coefficient, meaning that whether a squirrel kuks or not does not influence the value of the coefficient for `Above Ground Sighter Measurement`.

The model does not have a high predictive power as shown by the low AUC (0.512) and accuracy  $((113+179) / (113 + 124 + 145 + 179) = 0.52)$ . This is somewhat expected given



what we saw in the EDA section, where the distribution of the predictors is either similar across the two levels of response or very imbalanced.

## Discussion and Conclusion

In this paper, we have investigated which factors significantly affect whether a squirrel is indifferent to human presence, using a logistic regression model and a 10-fold cross-validation for model selection. We concluded that based on our logistic model, there are three significant predictors of **Indifferent**, which are **Kuks**, **Quaas**, and **Above Ground Sighter Measurement** since their p-values in the model is less than the  $\alpha = 0.05$  significance level. Putting this into a larger context, this shows that the squirrels' attitude to humans are largely caused by human presence since **Quaas** are sounds squirrels make when they see a ground threat (either a human or perhaps a pet dog), and our model shows that its presence is predicted to decrease the squirrels' odds of being indifferent. Surprisingly, both **Tail Flags**, which occurs when squirrels identify a threat and seek to confuse it, and **Tail Twitches**, which indicate curiosity, were not significant in the model. The insignificance of these variables may indicate that the squirrels do not see humans as a threat and are more preoccupied with other threats like hawks and dogs.

Although **Above Ground Sighter Measurement** has a statistically significant coefficient, the actual value of the coefficient is near-zero. The fact that it is negative is difficult to interpret because it suggests that the higher the squirrel is from the ground, the lower the odds of it being indifferent to the researchers. However, logically, we might argue that the higher the squirrel is off the ground, the safer it should feel from the researchers. Or perhaps, there are confounding factors: maybe the squirrels on trees are skittish due to another threat.

One important limitations is the concerns about the way data is collected. For example, based on the data description, the data is collected from volunteer sighting in NYC. This might introduce some lack of rigorousness, especially when it comes to numeric variables like **Above Ground Sighter Measurement** since the volunteers might give a very subjective estimation of the height above the ground when standing and observing from distance. In addition, there are a lot of NA's in the dataset, which to some extent limits our ability to do analysis.

For future work, we will try to extract more useful predictors from feature selection and feature engineering steps in order to increase the predictive power of our model. For example, instead of simply dropping the longitude and latitude from the dataset, we could potentially search the longitude and latitude of downtown NYC and use this information to calculate the Euclidean distance (or perhaps Manhattan distance since it's in NYC) between each observation and the city center for the model. In addition, we could try some non-linear machine learning models, such as random forest and boosting, since the linear model seem not to be able to give accurate predictions in this case.