

# Udacity Data Analyst Nanodegree

## P3: Wrangle and Analyze Data

*Author: Challa Sri Venkata Divya Madhuri*

*Date: Feb 4, 2018*

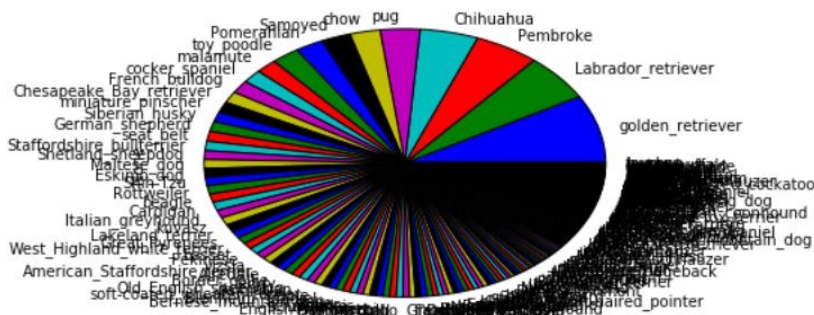
### Analysis and visualization

1. Plotting piechart to find the difference in 'Possible\_breed\_1' column between clean dataset (having only dog entries) and unclean dataset

To find the difference in distribution of the column 'Possible\_breed\_1' before and after cleaning, separate pie-charts are plotted for the 'Possible\_breed\_1' variable for both clean dataframe and the old dataframe before cleaning.

```
In [42]: import matplotlib.pyplot as plt
%matplotlib inline
# plot all possible dog breeds for df2_clean which has only dogs
df2_clean.Possible_breed_1.value_counts().plot(kind='pie')
```

<matplotlib.axes.\_subplots.AxesSubplot at 0xf3cc390>



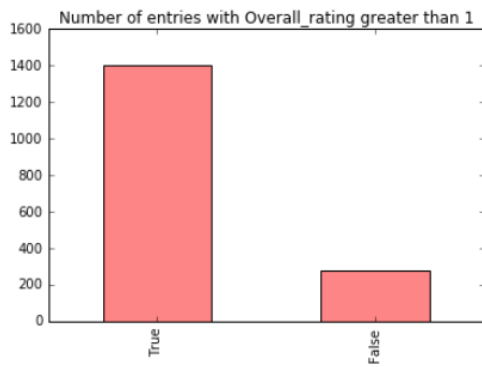
```
In [46]: # plot all possible dog breeds for df_unclean which has things other than dogs
df_unclean.Possible_breed_1.value_counts().plot(kind='pie')
```

**2. Analyze 'Overall\_rating' column to find number of entries with Overall\_rating >=1**

```
In [47]: # plotting number of entries having overall_rating >=1
df2_clean['rating_morethan_one'] = df2_clean.Overall_rating >=1
df2_clean.rating_morethan_one.value_counts()
```

[illegible]

```
<matplotlib.axes._subplots.AxesSubplot at 0x12849b00>
```



### 3. Insights from the Analysis

- As seen from the two piecharts plotted above, the clean dataset i.e df2\_clean has higher proportion of dog breeds like golden\_retriever, Labrador\_retriever, Chihuahua, Pembroke etc , but there is comparatively lesser proportion of dogs in df\_unclean (the unclean dataset).
- Only 280 entries of dogs have Overall\_ratings less than 1 , and the rest 1406 entries have Overall\_rating greater than 1.
- The 5 most popular dog breeds from the above dataset are golden\_retriever, Labrador\_retriever, Pembroke, Chihuahua and pug.

## Storing Data

The final cleaned dataframe df2\_clean is stored as CSV file named twitter\_archive\_master.csv

```
In [37]: df2_clean.to_csv('twitter_archive_master.csv')
```