Michael Albers
Project 5: Data Gathering
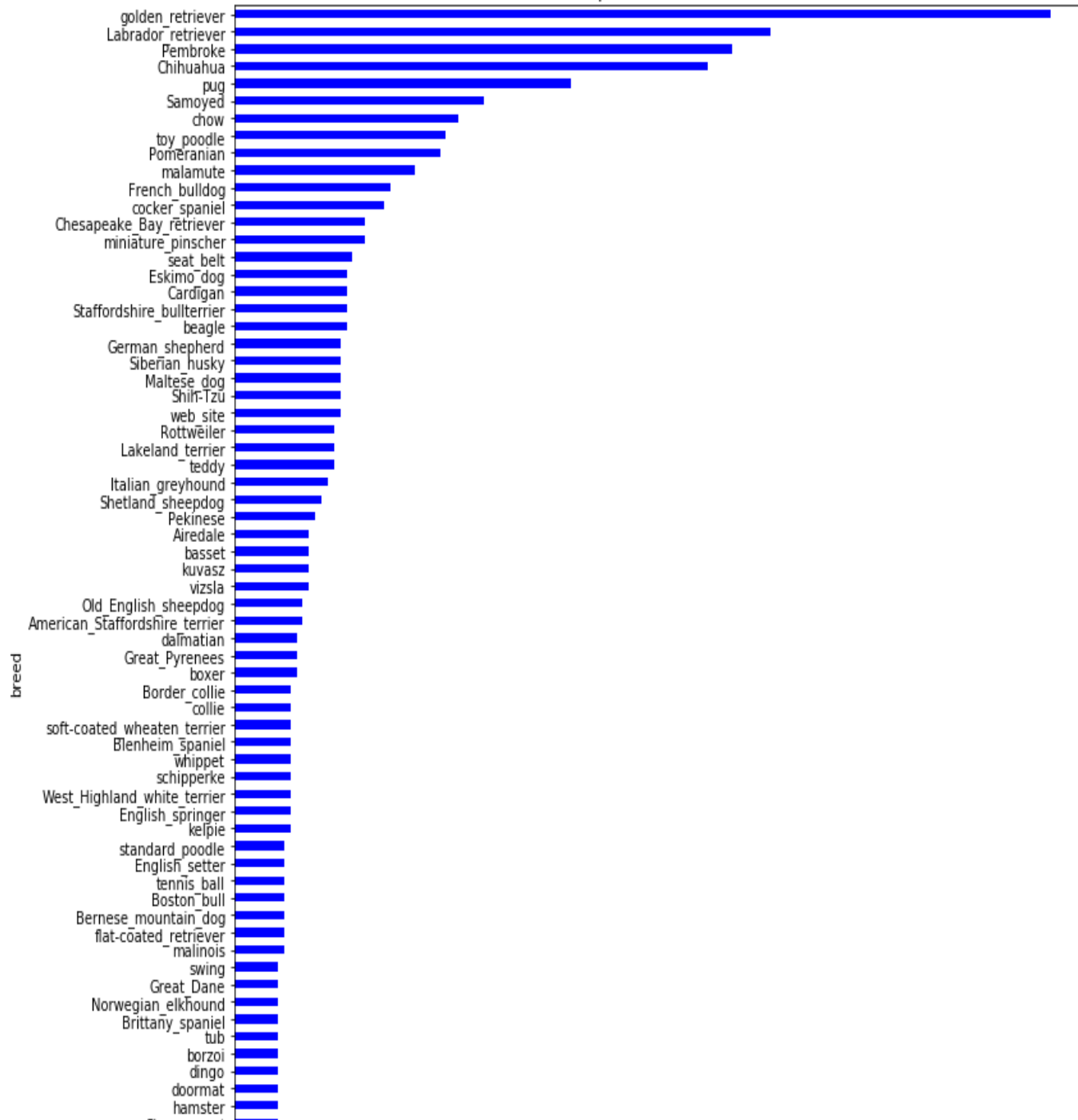
## Understanding Twitter "We Rate Dogs" breed prediction algorithm

Hey dog enthusiasts!  Twitter has an interesting site called We Rate Dogs https://twitter.com/dog_rates.  Dog lovers from all over can submit photos of their dogs along with a rating and brief description of their adorable dog.  An algorithm exists which takes submitted photos of people's dogs and makes three separate predictions of the dog's breed from one of the images submitted.  This blog post attempts to evaluate the accuracy of the neural network prediction algorithm for a tweet's dog.

The algorithm makes three predictions for each tweet's dog.  The first prediction usually has the highest confidence percentage of being right.  I was interested in seeing which predicted breed had the most "first" predictions.

As you can see from below bar chart, Labradors and Golden retriever breeds have the most "first" predictions.  But also notice that the algorithm makes a few "false" predictions.  For example, it predicts a few tweet posts as a "tub" or maybe a dog named "tub" for a first prediction.

First prediction count > 20

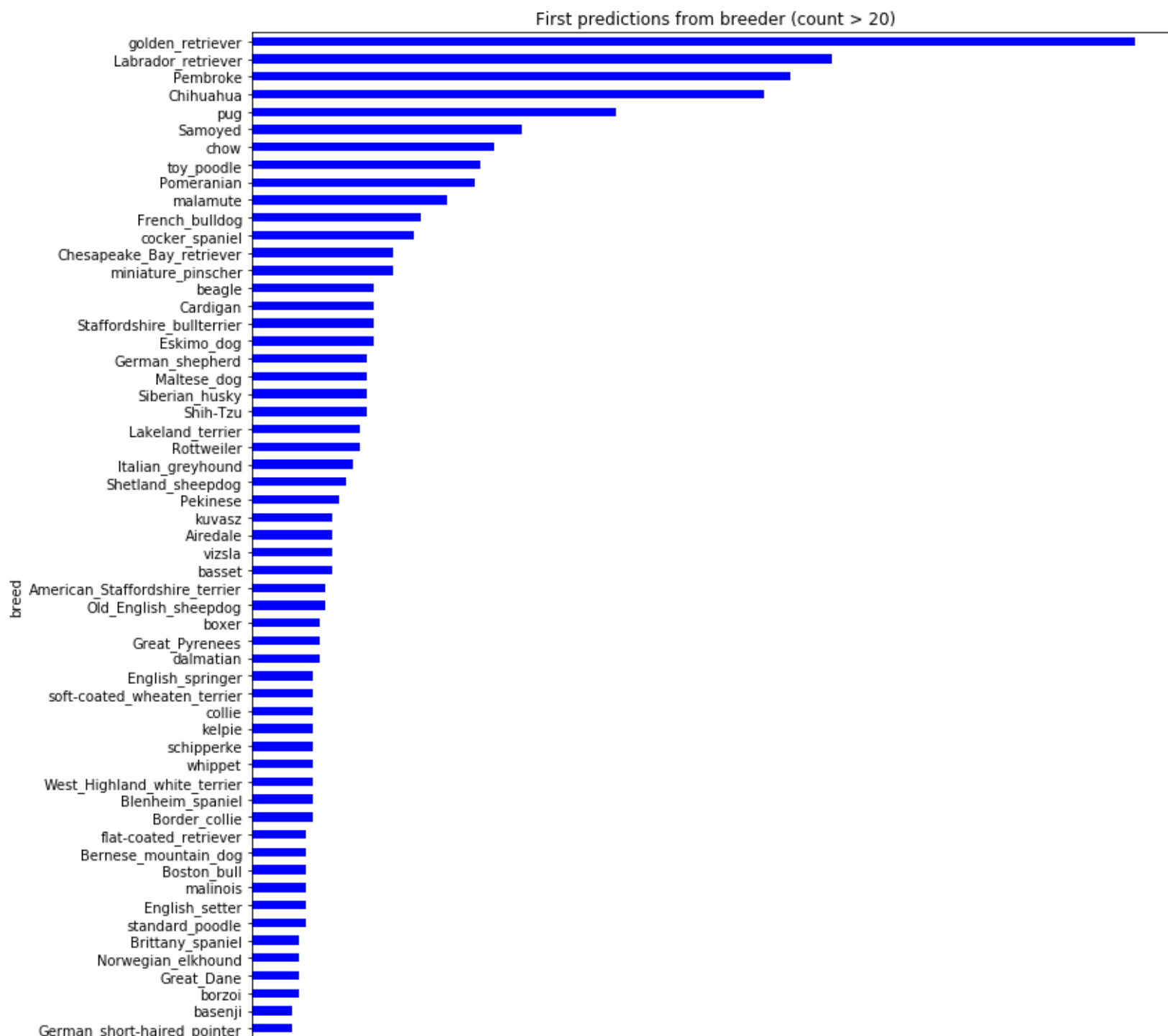| breed | |
|---|---|
| golden_retriever | |
| Labrador_retriever | |
| Pembroke | |
| Chihuahua | |
| pug | |
| Samoyed | |
| chow | |
| toy_poodle | |
| Pomeranian | |
| malamute | |
| French_bulldog | |
| cocker_spaniel | |
| Chesapeake_Bay_retriever | |
| miniature_pinscher | |
| seat_belt | |
| Eskimo_dog | |
| Cardigan | |
| Staffordshire_bullterrier | |
| beagle | |
| German_shepherd | |
| Siberian_husky | |
| Maltese_dog | |
| Shih-Tzu | |
| web_site | |
| Rottweiler | |
| Lakeland_terrier | |
| teddy | |
| Italian_greyhound | |
| Shetland_sheepdog | |
| Pekinese | |
| Airedale | |
| basset | |
| kuvasz | |
| vizsla | |
| Old_English_sheepdog | |
| American_Staffordshire_terrier | |
| dalmatian | |
| Great_Pyrenees | |
| boxer | |
| Border_collie | |
| collie | |
| soft-coated_wheaten_terrier | |
| Blenheim_spaniel | |
| whippet | |
| schipperke | |
| West_Highland_white_terrier | |
| English_springer | |
| kelpie | |
| standard_poodle | |
| English_setter | |
| tennis_ball | |
| Boston_bull | |
| Bernese_mountain_dog | |
| flat-coated_retriever | |
| malinois | |
| swing | |
| Great_Dane | |
| Norwegian_elkhound | |
| Brittany_spaniel | |
| tub | |
| borzoi | |
| dingo | |
| doormat | |
| hamster | |

Without knowing the neural algorithm and how it works, it is hard to understand why these 'false' predictions appear.  The next bar chart is limited to the number of first predictions being less than 20 for a predicted breed.  As you can see in bar chart, below there are far more "false" breed predictions than before.  However, note that the scale for the second bar chart goes up to about 18 compared to a 550 in the chart above.  So even though there are false predictions of breeds, the number is quite small comparted to the roughly 1,000 or more "first" predictions in our dataset.
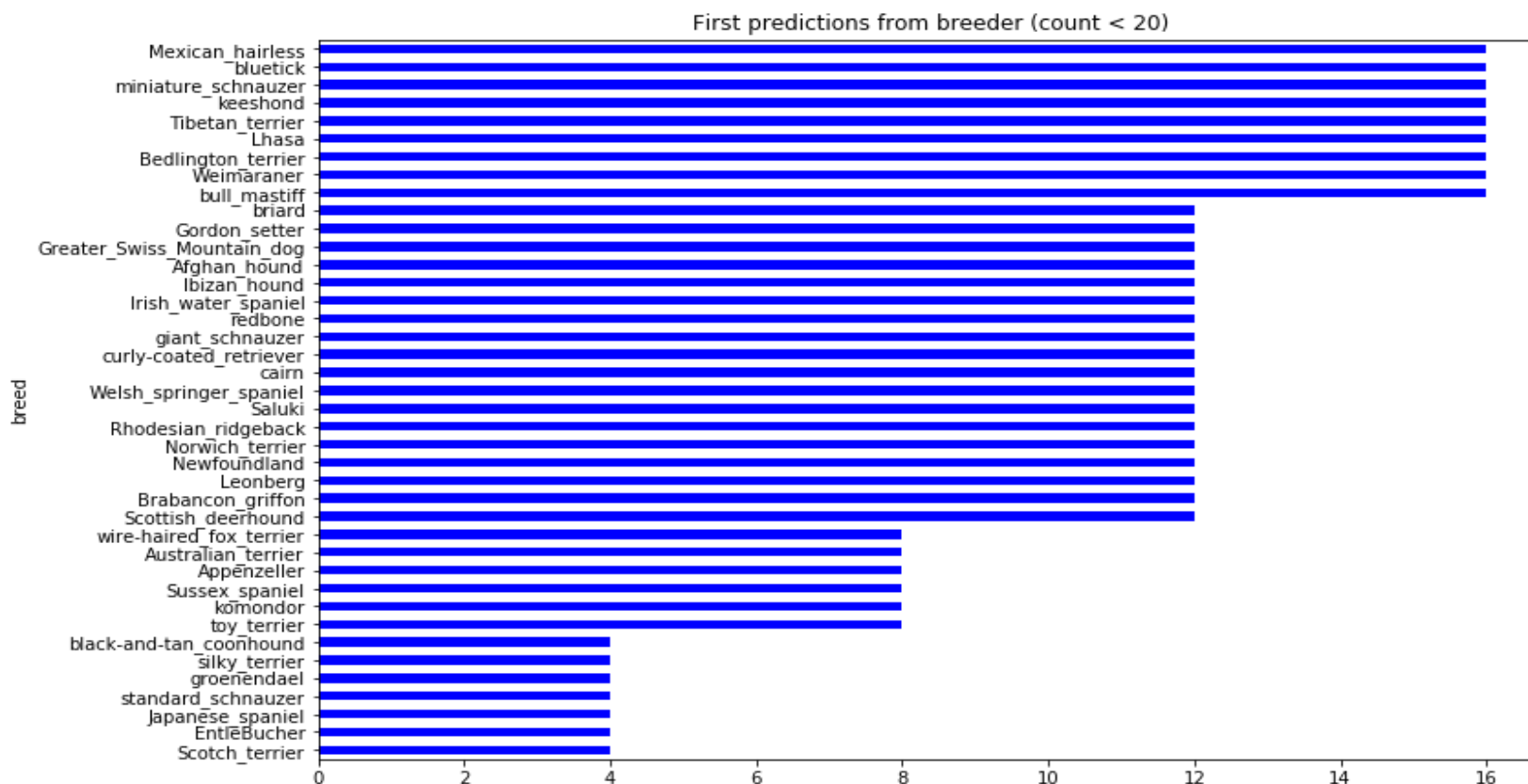
First prediction count < 20

bull_mastiff
bathtub
home_theater
hog
shopping_cart
miniature_schnauzer
Mexican_hairless
Tibetan_terrier
Weimaraner
barrow
Lhasa
keeshond
jigsaw_puzzle
Arctic_fox
Bedlington_terrier
goose
brown_bear
guinea_pig
bluetick
vacuum
washbasin
refrigerator
motor_scooter
Afghan_hound
triceratops
ox
sea_lion
prison
space_heater
ski_mask
white_wolf
redbone
window_shade
stone_wall
mousetrap
wood_rabbit
wombat
seashore
ram
comic_book
common_iguana
Scottish_deerhound
cairn
Irish_water_spaniel
hippopotamus
giant_schnauzer
curly-coated_retriever
cowboy_hat
Arabian_camel
Christmas_stocking
Gordon_setter
Greater_Swiss_Mountain_dog
Ibizan_hound
Brabancon_griffon
Leonberg
Norwich_terrier
Rhodesian_ridgeback
Saluki
Welsh_springer_spaniel
balloon
bow_tie
briard
Newfoundland
laptop
koala
komondor
wire-haired_fox_terrier
Appenzeller
paper_towel
birdhouse
leatherback_turtle
snorkel
jellyfish
Australian_terrier
hen
jack-o'-lantern
patio
sorrel

We want to find a reason or a condition as to why these extra "false" predictions appear on our second bar chart. While there are not many, it does create extra "noise" on our bar chart and may give a false impression that the neural prediction algorithm is faulty.  The next chart, brings in an additional variable called "is bred" which indicates whether the dog in the twitter post comes from a breeder. Here we see, there are no "false" first predictions with the additional flag "is bred".
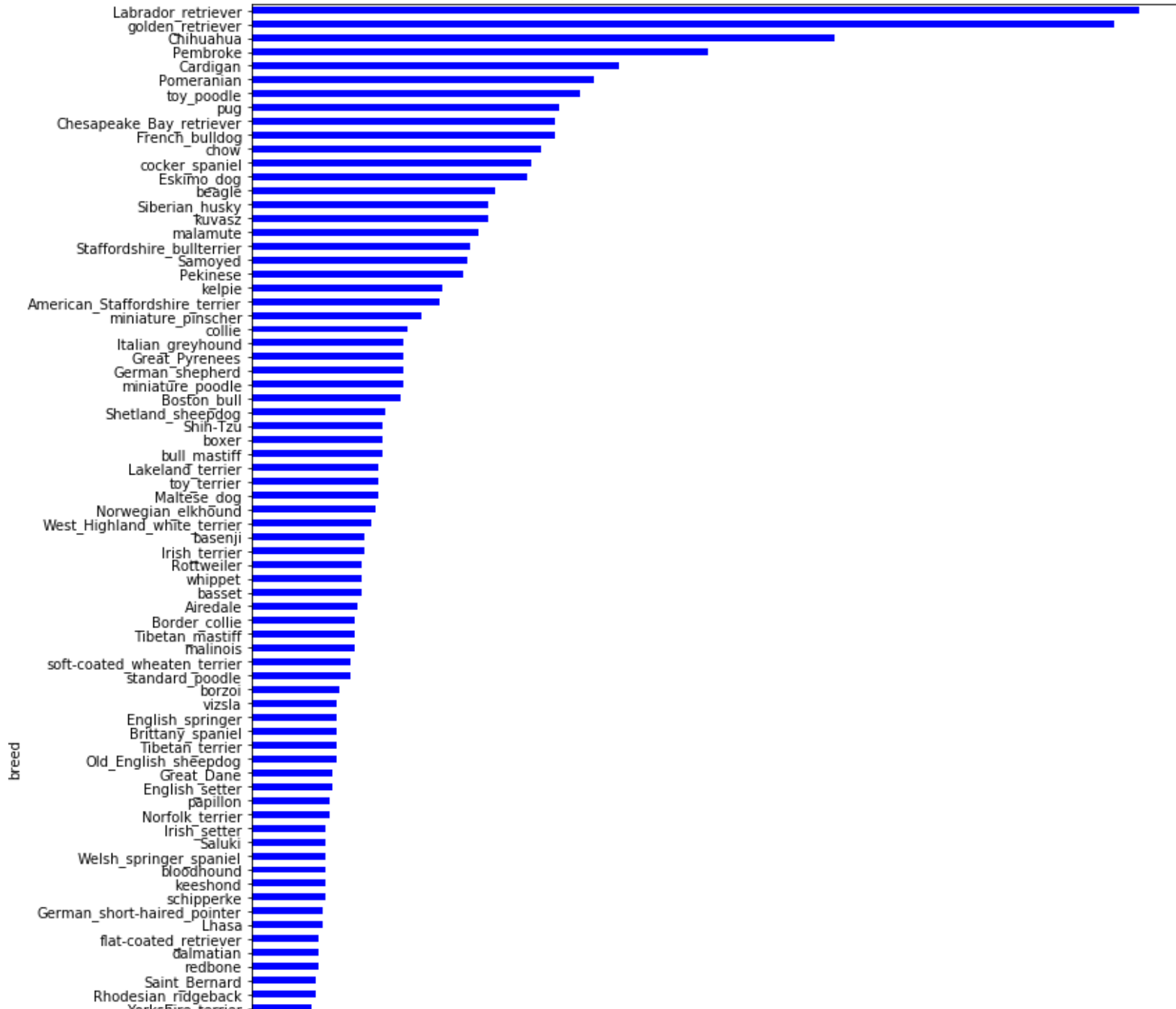
First predictions from breeder (count > 20)

Let's now view cases where the number of "first" predictions of a breed is less than 20 with the additional "is bred" flag. Here again, we see a big improvement in that there are no "false" first predictions.  It appears the additional "is bred" flag is making our charts look much more accurate for breed predictions.

First predictions from breeder (count < 20)



However, thus far we have looked only at "first" predictions.  Recall, that each tweet post for a dog gets three predictions.  Our next chart will chart predictions based on "is bred" being true for all predictions (first, second,
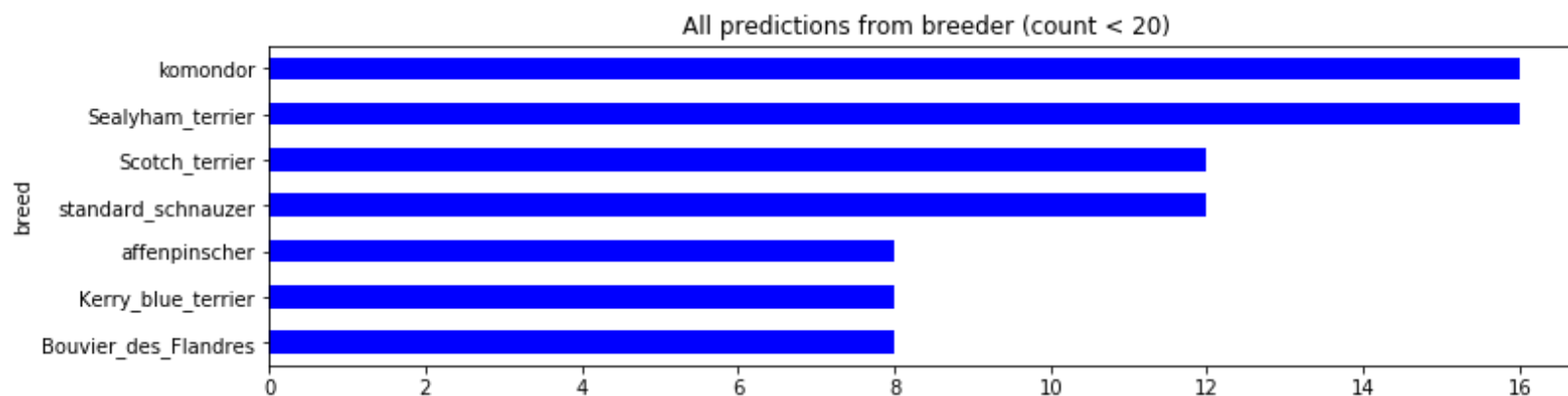
and third).  We will try to determine if the neural network algorithm for predicting a dog's breed is better for twitter posts in which the "is bred" flag is True.  Again, we do not see any false predictions across all three predictions for twitter posts of dogs from a breeder.

All predictions from breeder (count > 20)

Lastly, the next chart displays predictions of breeds where the number of overall predictions is less than 20.

Again, we don't see any false predictions for twitter posts with dogs a from a breeder.



All predictions from breeder (count < 20)

## Conclusion

We see that the predicted breeds are very accurate when we limit our data to just dogs from a breeder. This limits the effectiveness of the prediction algorithm to just bred dogs. There are many dogs that are not from a breeder. It is worth nothing again that the number of false predictions were small but it appeared to be many because our bar charts showed many different predictions of erroneous things like "web_site", "dishwasher". This extra noise may scare someone away from trusting this algorithm. The final question is why does the algorithm work so much better with "is bred" dogs? Is it possible, the neural network algorithm for predicting breeds is somehow linked to a breeder database network of images which allows it to match submitted images with a breeder database? This is only a guess on my part as we do not know the algorithm was created.