1. Compare the ratings of different dog stages—Which dos stage has the highest rating? Are the differences statistically significant?

First, I applied "groupby" function to calculate the ratings of different dog stages, then plotted the findings into a bar chart.

Among all the dog stages, "doggo, floofer" has the highest average rating. However, there are only two such cases, and their ratings are both 13, so it's not really suitable for statistical analysis. Therefore, I chose the second highest stage--puppo. By comparing it with different stages, the conclusion can be drawn that puppo has higher ratings than pupper and unspecified dogs. It is not significantly higher than the ratings of doggo and floofer.

Then I compared the ratings of the ones with specified dog stages and the ones with specified dog stages. I applied hypothesis testing, with the null hypothesis being that the two groups have the same rating. As a result, the group with specified dog stages has higher ratings.

2. Compare the numbers of favorites and retweets of different breeds.

I applied the "groupby" function to calculate the average number of favorites and retweets of different breeds, and then plotted my findings into two bar charts. From the charts, we can draw the conclusion that toy poodle is the most popular breed.

3. Is rating correlated with the numbers of favorites and retweets?

I applied the "groupby" function to calculate the average retweets and favorites of different ratings. In general, there is a positive correlation with rating and the numbers of favorites and retweets. Then, I applied linear regression to examine whether the positive correlation is statistically significant. Both linear regressions show that rating has a positive impact on the numbers of retweets and favorites.