

### Preprocessing

I use pandas to handle the CSV by simply reading it in. I then turn the resulting dataframe to a numpy array, then define some variables to be used.

### Question 1

Here, we compute the Fleiss' Kappa inter-annotator agreement. I first use a for loop to populate the Fleiss matrix as shown [here](#) in the table under Worked Example. The process was fairly simple. First, I initialized a  $N \times k$  matrix of zeros where  $N$  is the number of subjects (2199 subjects) and  $k$  is the number of categories (3 categories). I then looped through the Fleiss Matrix and filled it by adding a one to the categories by each rater. From there, I simply used statsmodel's fleiss\_kappa function with my fleiss matrix as the argument. The result was a Fleiss' Kappa score of 0.650, meaning the ratings have "substantial agreement" based on Table 1 in the assignment

### Question 2

Here, we compute Cohen's Kappa between each annotator and the majority vote. To calculate Cohen's Kappa, I simply looped through the ratings and ran sklearn's cohen\_kappa\_score(ratings1, rating2) function with ratings1 being the annotator's rating and rating2 being the majority votes. The scores are as follows, ranked in order of their agreement with their Cohen's Kappa value.

Annotator	Cohen's Kappa
1	0.895
3	0.798
2	0.784

Note, I also computed the Cohen's Kappa between all combinations of the annotators. I left it because I found the result to be fairly interesting. It was interesting to note that Annotator 1 had the highest two highest Cohen Kappa in this experiment, since it also had the highest Cohen's Kappa score when compared to the majority vote.

### Question 3

Here, we recalculate the majority vote by giving Annotator 1's vote to be weighed twice as much as Annotators 2 and 3. The premise is simple, treat Annotator 1's vote as two votes, meaning that the majority vote is calculated from 4 votes rather than 3 votes. To calculate the weighted majority vote, I loop through the subjects, and calculate the new weighted majority votes for each subject. There are two conditions in the for loop: (1) when rating2 equals rating3, but rating1 does not equal rating2 (and rating 3 since they are equal) (2) all other cases. Condition (1) describes a scenario of a tie in the weighted majority vote. To deal with ties, I chose to

randomly choose between rating 1 and the other rating (i.e. rating 2 and rating 3). This was to balance the expertise of rater 1 and the actual majority of the annotators. The other condition is when there is no tie in the weighted majority vote. The majority vote won't actually change in this scenario, but I chose to do it to make it easier. From here, I write to a csv file using numpy savetxt, which saves things into a CSV file.

#### Question 4

Two cues that I observed that may have correlations with sentiment are words with sentiment and mouth configurations. These cues carry a loose interpretation, but words with sentiment refer to happy, sad, bad, good, etc. and mouth configuration refers to smiles, frowns, mouth resting position. Below is how I encoded them.

	Positive (+1)	Neutral (0)	Negative (-1)
Sentiment words	Words like good, bad, etc.	Description words of the movie, words like average.	Words like bad, awful, boring, etc.
Mouth configuration	Smile	Resting position	Frown

I labeled the videos according to the above table. For example, if the video contains an utterance like "this movie was amazing!" I would label the sentiment word with a 1 corresponding to a positive word. I hypothesize that both my cues have positive correlation with sentiment indicated by the majority vote.

#### Question 5

Here, I label the clips for 0h-zjBukYpk \* and 1DmNV9C1hbY \* myself based on the behaviors outlined in the previous step. I used pearson(dataX,dataY) in Excel to calculate the Pearson Correlation Coefficient. I found the Pearson Correlation Coefficient to be 0.753 for sentiment words and .519 for mouth configuration. Both these coefficients indicate large strength of association according to [this article](#). I believe that the coefficient for sentiment words was greater than the coefficient for mouth configuration because mouth configuration can either be misleading or not present. For example, suppose one is making fun of the movie using words like "this movie was so so so bad", but they are smiling because they found it amusing how bad the movie was. This is likely to be rated as negative sentiment overall. Below are my calculations

			Pearson Coefficients	
Sentiment Word	Mouth Configuration	Majorith Vote	Sentiment Words	Mouth Configuration
0	0	0	0.7532165572	0.5188956412
-1	-1	-1		
0	0	0		
0	-1	-1		
0	0	0		
0	0	0		
1	0	1		
0	0	0		
0	0	0		
0	-1	0		
0	0	0		
0	0	0		
0	0	0		
1	0	0		
0	0	0		
1	0	0		
0	0	0		
-1	0	-1		
0	-1	0		
0	0	0		
1	1	1		
1	1	1		
1	1	1		
1	1	1		
1	-1	1		
-1	0	-1		
1	0	1		
0	0	0		
1	1	1		
-1	0	0		
-1	0	-1		

0	-1	-1		
1	0	-1		
1	0	1		
1	0	1		
-1	-1	-1		
0	0	0		
0	-1	0		
0	0	-1		