

1 Introduction

Team name: Beavers

Names: Brendan Flaherty, Sujit Iyer, Sayuj Choudhari

Work Division: Worked together on all parts of the project.

Packages Used: pandas, numpy, matplotlib, seaborn, counter, surprise

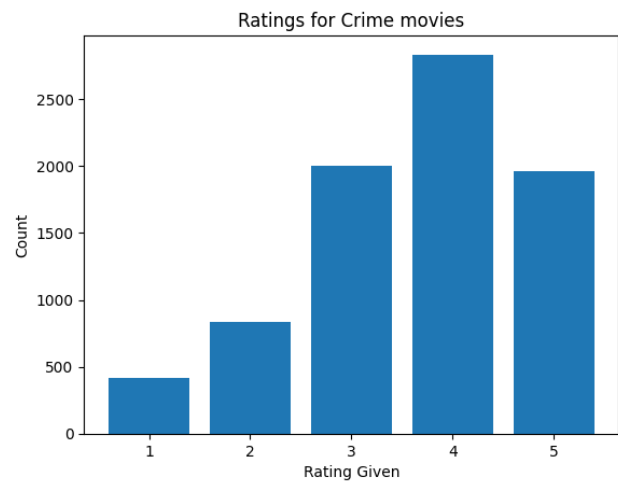
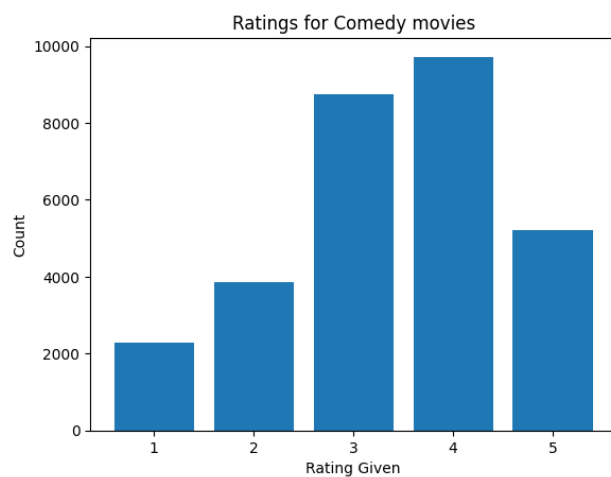
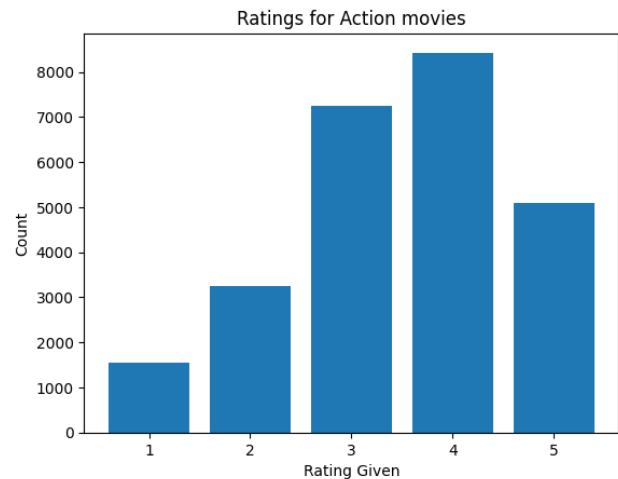
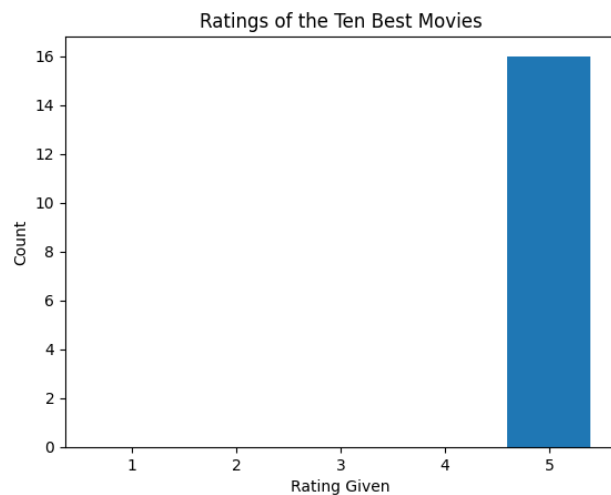
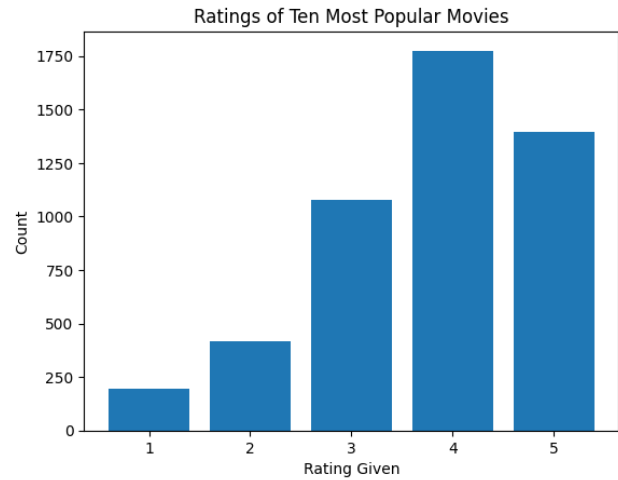
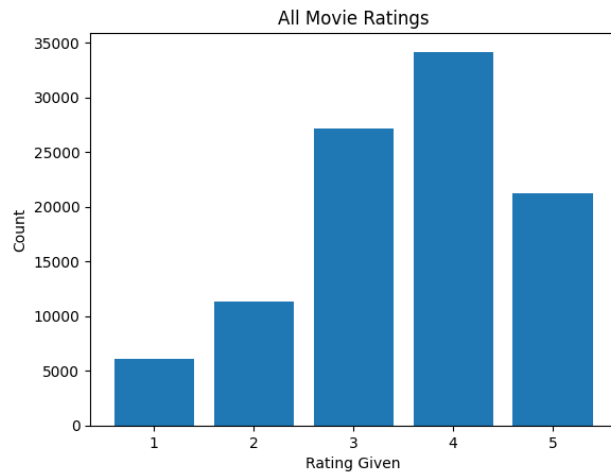
Google Colab Link: [beavers_mp2.ipynb](#)

Piazza Link: [Beavers' Piazza Post](#)

2 Basic Visualizations

We created visualizations for all ratings in the MovieLens Dataset, for the ten most popular movies, for the ten best movies, and all ratings from Action, Comedy, and Crime genres. First, the histogram of all ratings had a mode rating of 4 and the frequencies were skewed right – favoring higher ratings. We expected to see this skewed distribution due to positivity bias and our belief that people would be more likely to give very good ratings compared to very bad ratings. The ratings of the most popular movies followed a similar distribution pattern of the overall ratings as well. One difference was a higher frequency of rating 5, but the mode was still 4. These results matched our expectations because logically the most popular movies make up a large subset of the entirety of movie ratings and therefore the top 10 most popular movies distribution should roughly follow the distribution for all movies. We also plotted the ratings of the ten best movies and found that each of the 16 ratings was a 5 (the frequency of all other ratings was 0). This matched our expectations because we never established a lower bound on the number of user ratings for each movie, which is why some movies have a high mean rating despite only 1-2 reviews. The distributions of the ratings for the best movies and the most popular movies were very different.

We plotted the ratings of Action, Comedy, and Crime genres. The distributions of the ratings for all the genres were all very similar. For all genres, the rating 4 was the most common. Comedy and Action genres both had many more ratings than the Crime genre. There may not be lots of variation between the genre distributions since the genres are all very popular and enjoyed by many. The six graphs are included below:



3 Matrix Factorization Visualizations

Method choices and derivations

We created visualizations for 3 different matrix factorization methods: a self-produced SVD trained on the regularized loss function, a self-produced SVD with Bias variables for each movie and user trained on the regularized loss function, and Sci-kit Learn's surprise library functions for SVD. The first two SVD's were trained using stochastic gradient descent on the loss functions. The first method used the code from project 5, as we could apply the same loss function and gradient descent method to the SVD of the MovieLens dataset just with different settings of matrix and latent factor sizes. However, for the SVD with bias, the loss function changes to incorporate the bias terms at each element for user i and movie j , as the estimation for each matrix entry at i, j is given by $u_i^T v_j + a_i + b_j$ where a_i is user i 's bias and b_j is movie j 's bias. This results in new loss function:

$$\arg \min_{U,V} \frac{\lambda}{2} (\|U\|_F^2 + \|V\|_F^2) + \frac{1}{2} \sum_{i,j} (y_{ij} - u_i^T v_j - a_i - b_j)^2$$

Given the new loss function, the gradient with respect to each of the variables now changes and we must now take a gradient step for each of the biases as well. For the non-bias variables, notice that the chain rule is not affected by the addition of either bias term as both bias terms are added to the chain term and not directly multiplied or applied to any of the non-bias variables. Thus, we can make simple edit of just adding the bias terms to the chain term of the original SVD gradient:

$$\begin{aligned} \partial_{u_i} &= \lambda u_i + \sum_j (-v_j)(y_{ij} - u_i^T v_j - a_i - b_j) \\ \partial_{v_j} &= \lambda v_j + \sum_i (-u_i)(y_{ij} - u_i^T v_j - a_i - b_j) \end{aligned}$$

For the gradients with respect to the bias terms, notice that a_i, b_j are lone-standing terms of degree 1 with a -1 coefficient, thus it is a simpler chain rule and additionally the regularization term is removed as it has no dependency on the bias terms so it goes to 0 when differentiating, the resulting gradient gives:

$$\begin{aligned} \partial_{a_i} &= - \sum_j (y_{ij} - u_i^T v_j - a_i - b_j) \\ \partial_{b_j} &= - \sum_i (y_{ij} - u_i^T v_j - a_i - b_j) \end{aligned}$$

Notice for the code of the gradient step with stochastic gradient descent only takes steps with respect to a single point y_{ij} so the gradient computation is essentially the same for a_i, b_j so they share the same gradient descent step. This makes the computation a bit easier as we have 1 less gradient to compute but lessens the degrees of freedom the model has to train on as both biases are moving in the same direction.

The third method, Sci-kit Learn's surprise library was an off-the-shelf implementation which was very easy to implement given the ample online documentation. Like the majority of Sci-kit Learn's libraries, the

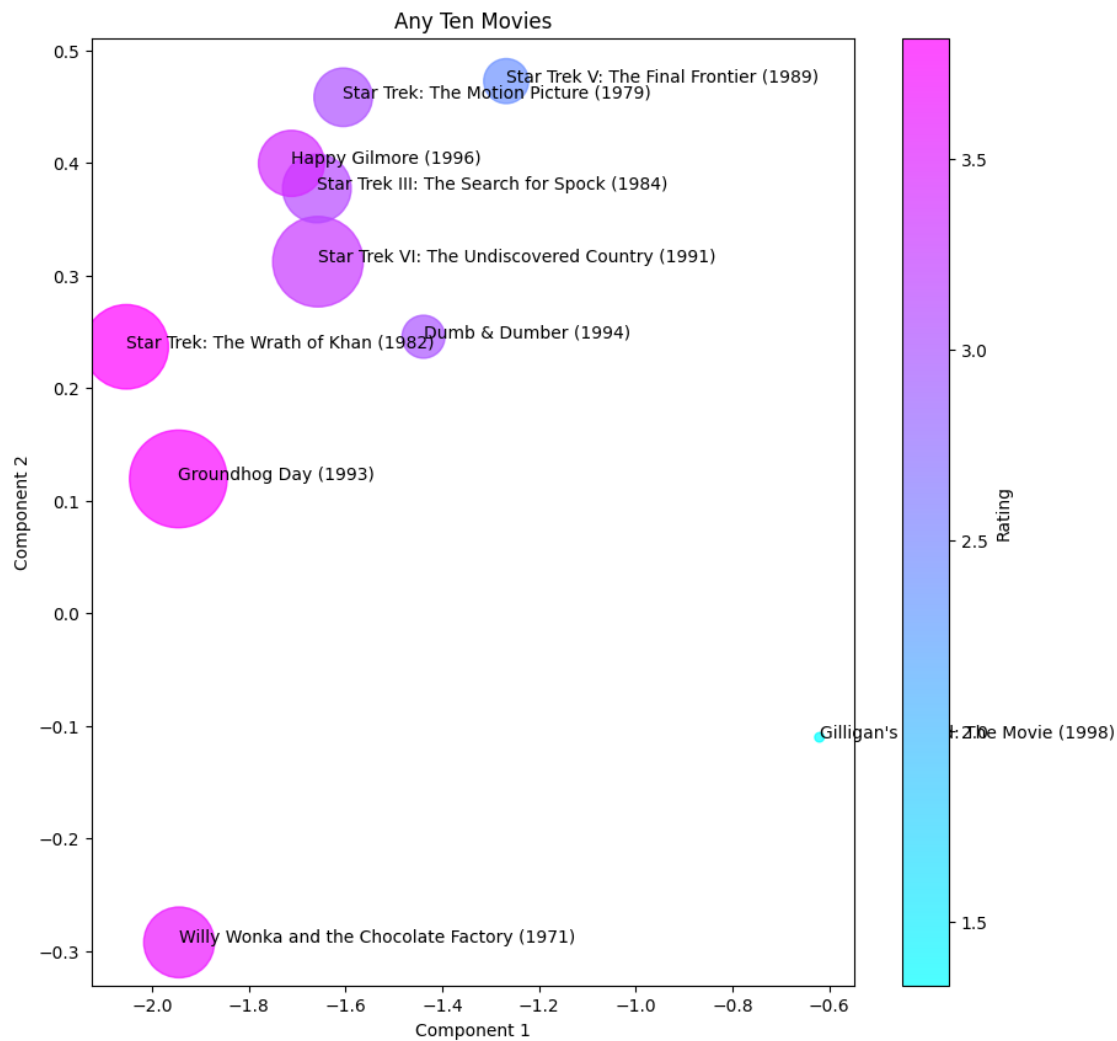
process was straightforward: Initialize an SVD object with desired hyperparameters, fit to the training data, and then find RMSE error by comparing the predictions on the test set. We also performed cross-validation on the entire data set to familiarize ourselves with the SVD library and compare the error to the other two methods.

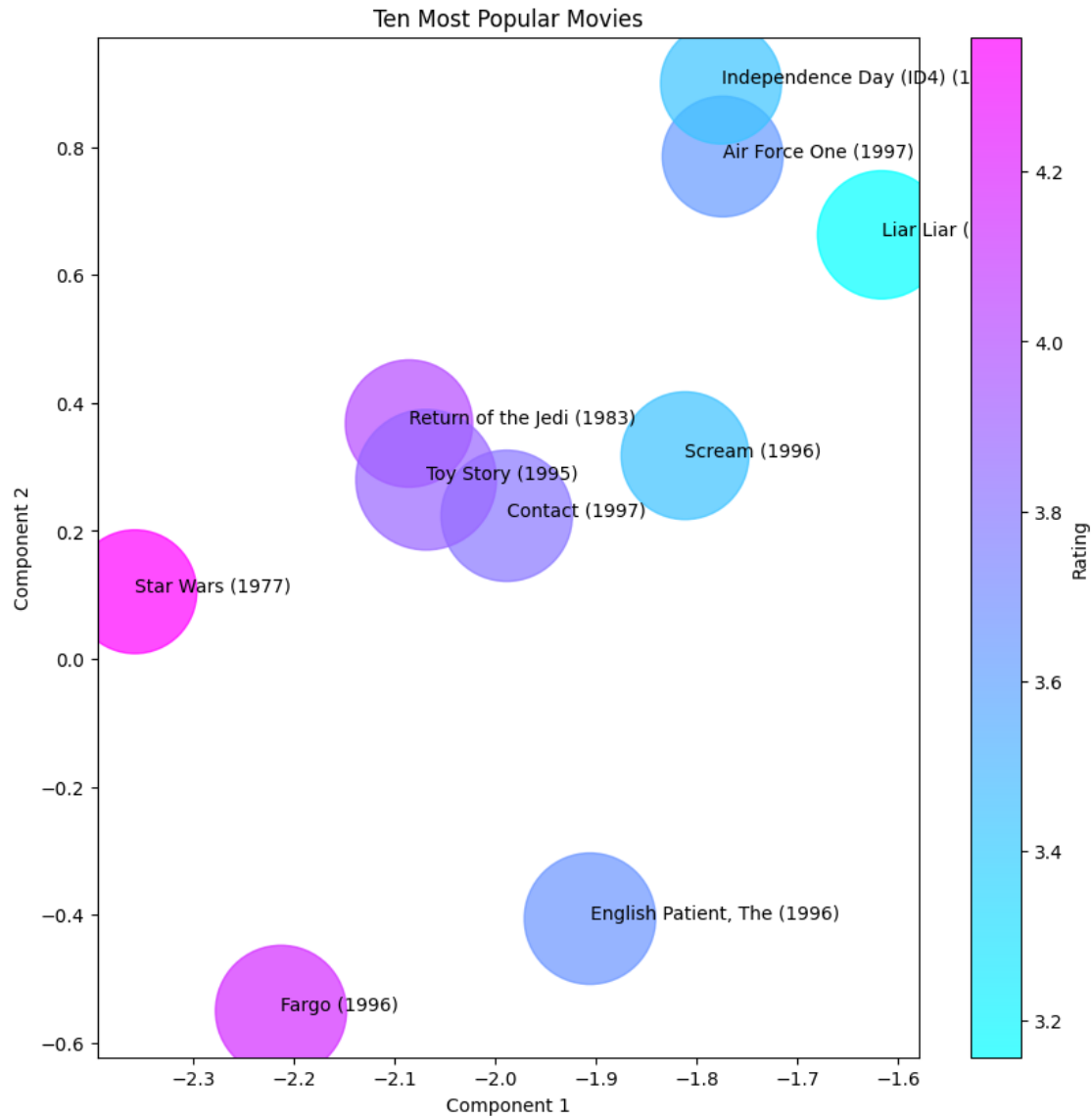
Generally, we saw that the test error was minimized for method 1 compared to methods 2 and 3. The most likely explanation is that methods 2 and 3 overfit the data. Methods 2 and 3 both had bias terms; adding bias terms adds a degree of freedom to the models, and therefore model complexity increases in expense to test accuracy.

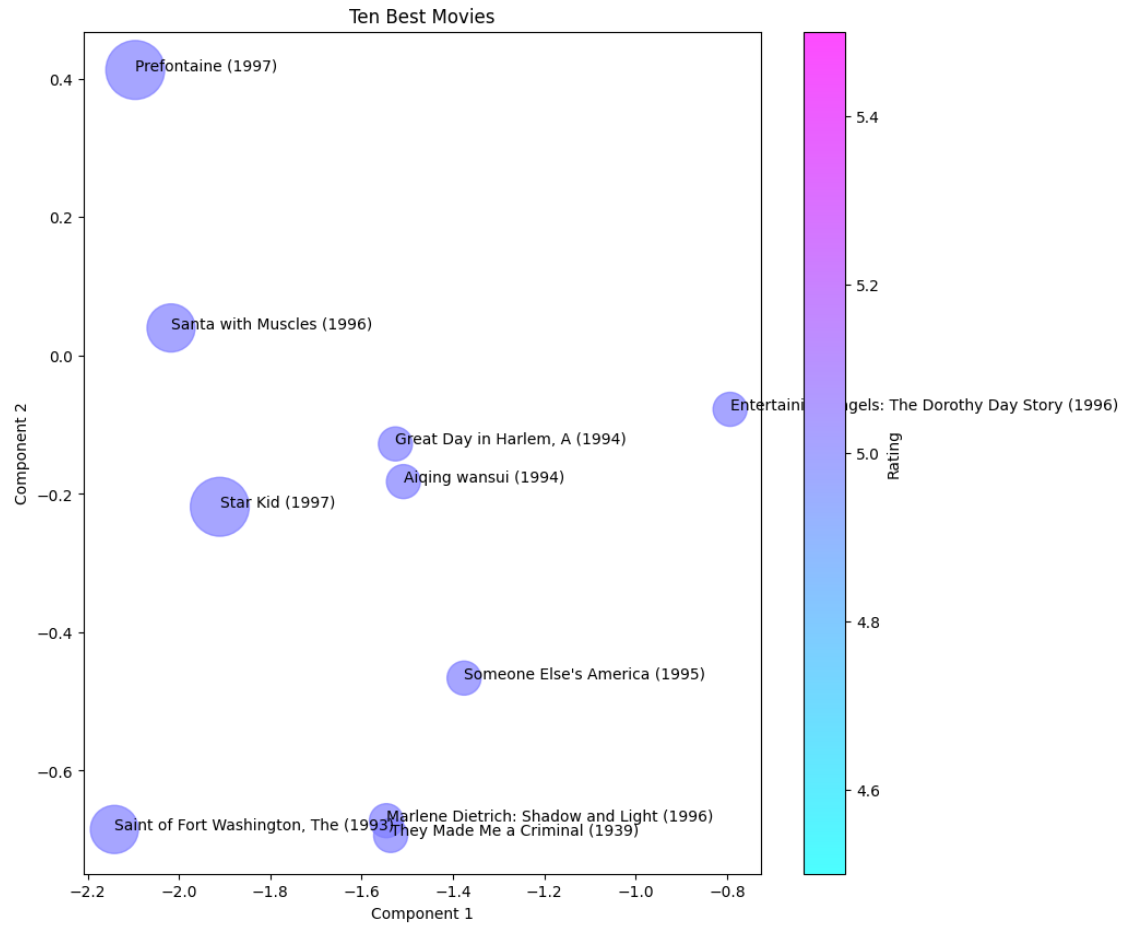
Using the methods described above, we produced 6 different plots for each method to compare the methods against each other. These visualizations plot given movies data projected against the first 2 latent factors of the SVD trained on the dataset. 10 movies were plotted on each graph and the choices of the 10 movies for each of the 6 graph types were: 5 star-trek movies and 5 comedy movies (to see if the latent factors clustered and split based off similarity), top 10 most popular movies (most ratings), top 10 best movies (highest average rating), 10 random action movies, 10 random comedy movies, 10 random crime movies. The plots are inserted below.

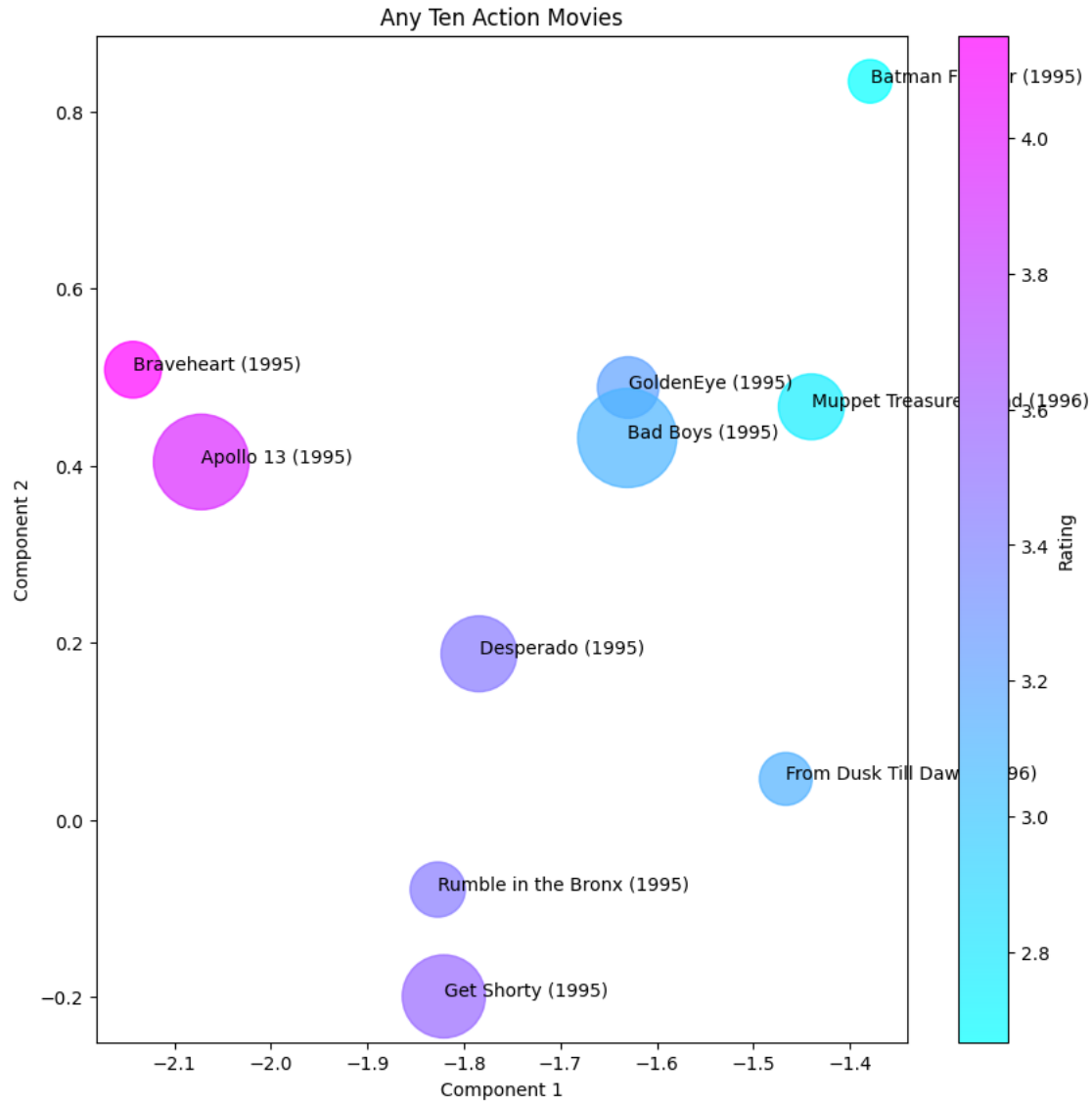
Method 1 Plots

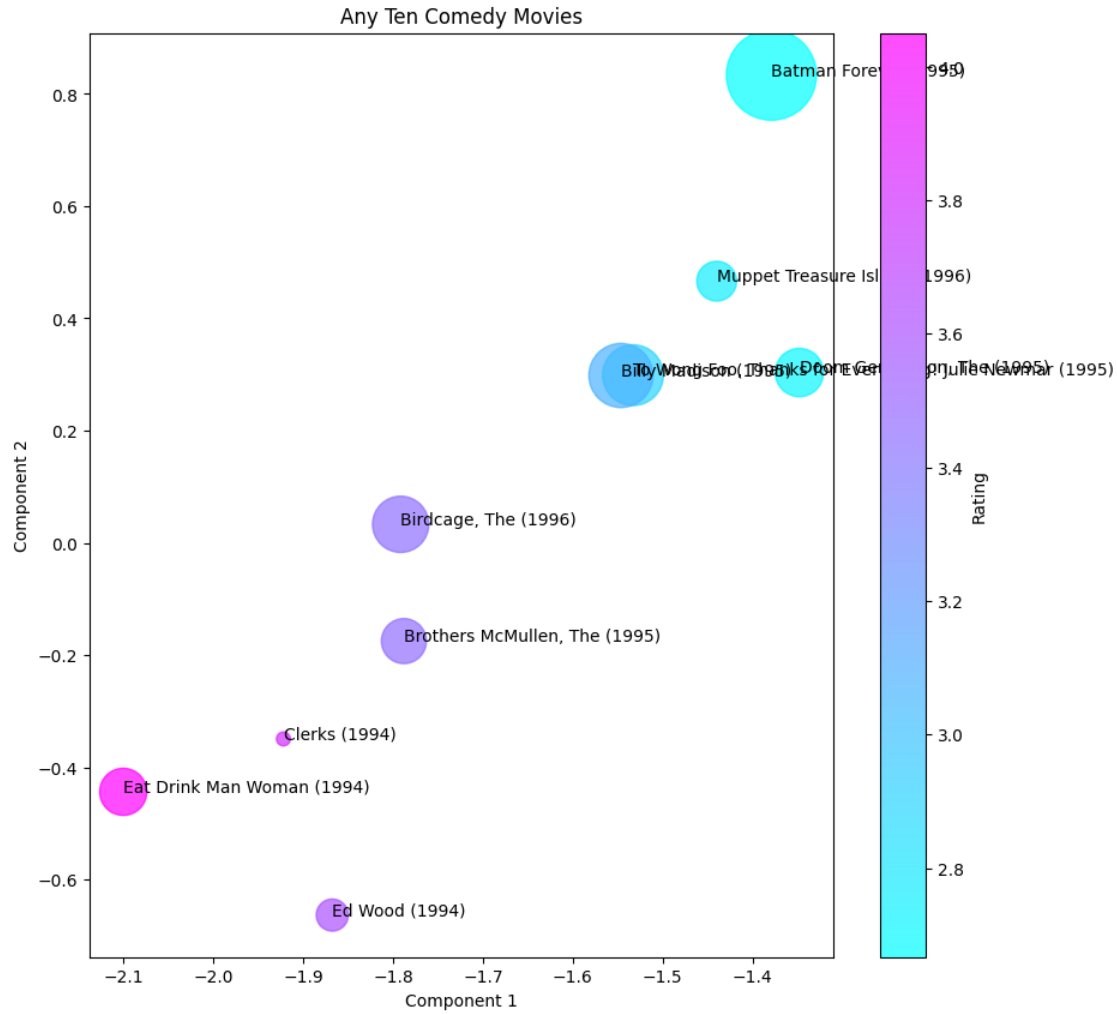
Below are the six visualizations resulting from Method 1.

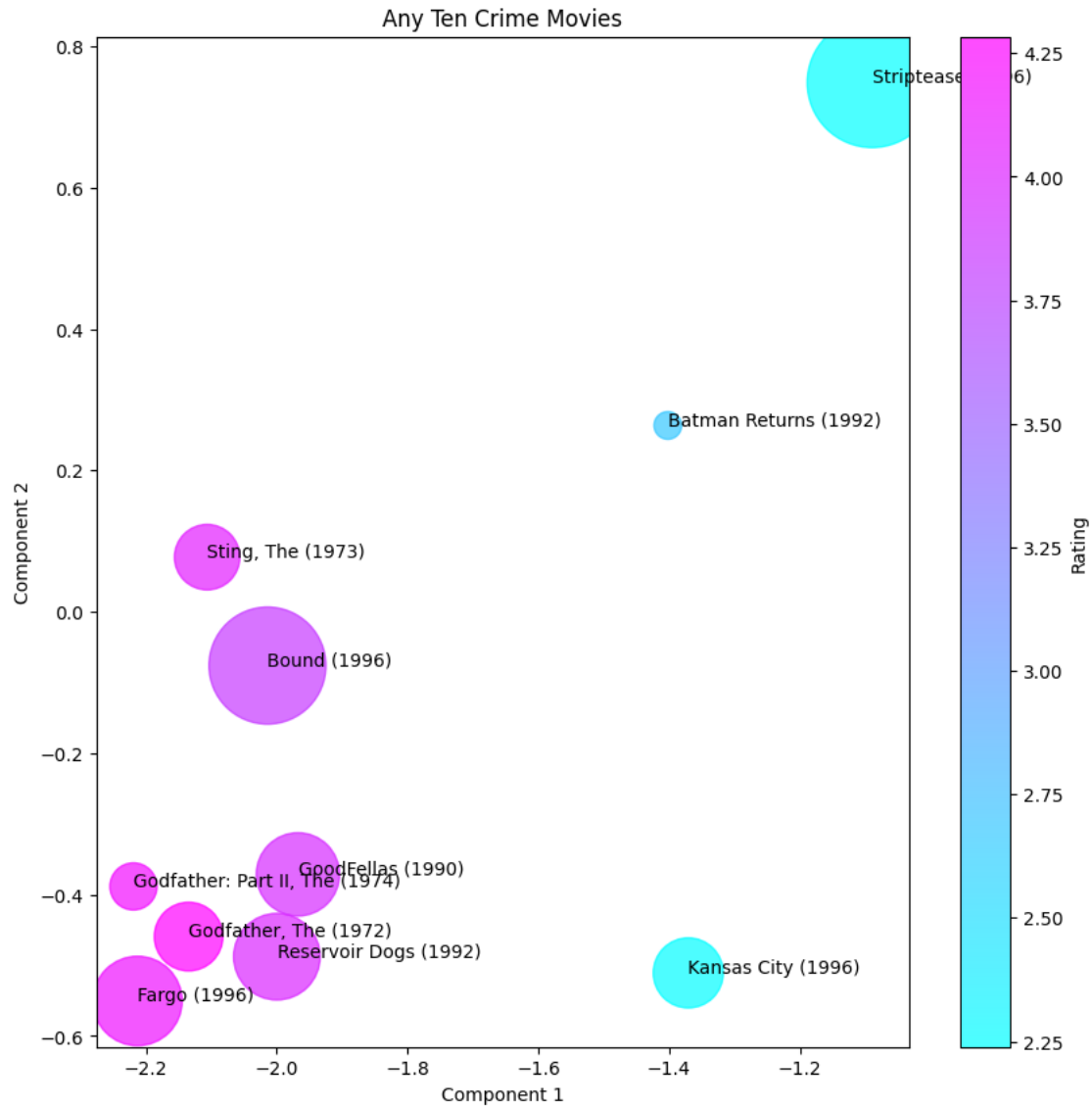






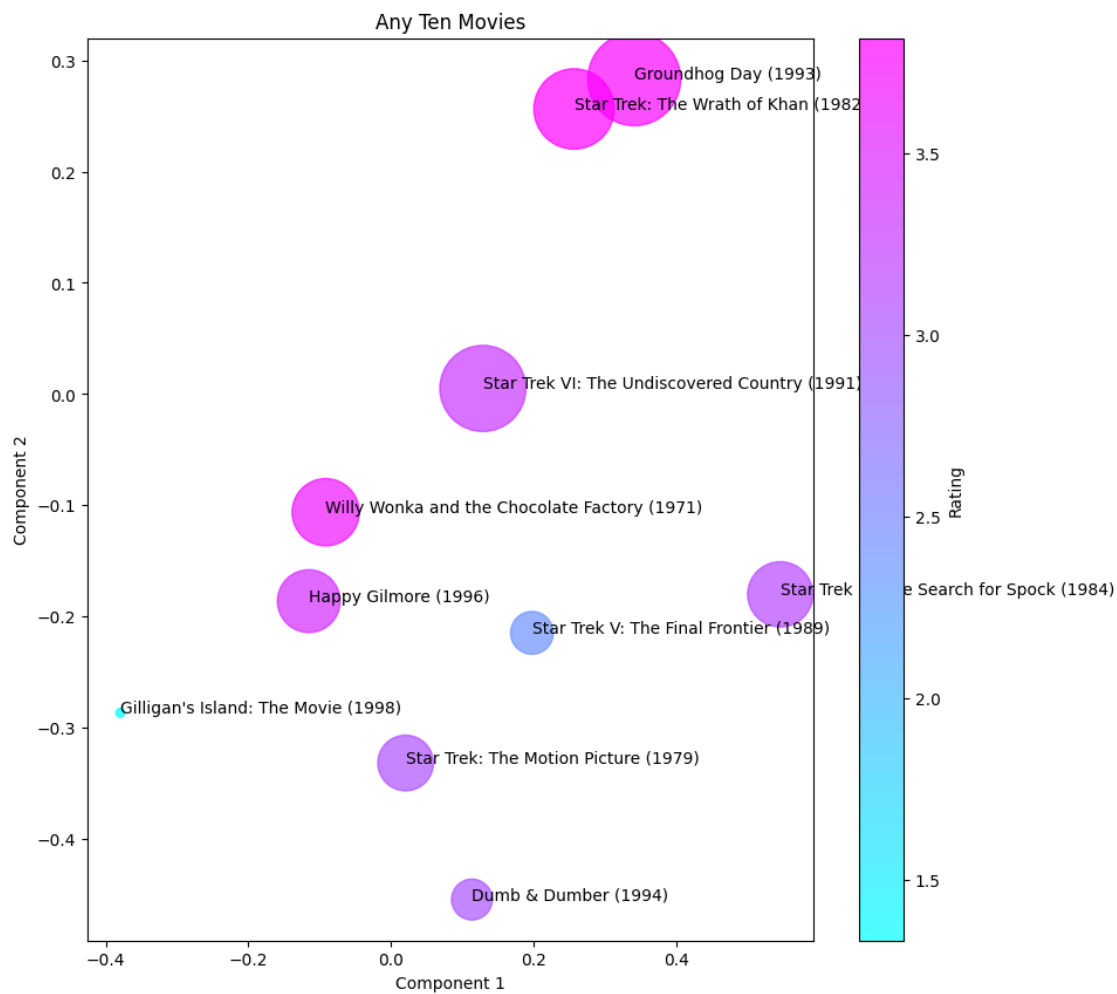


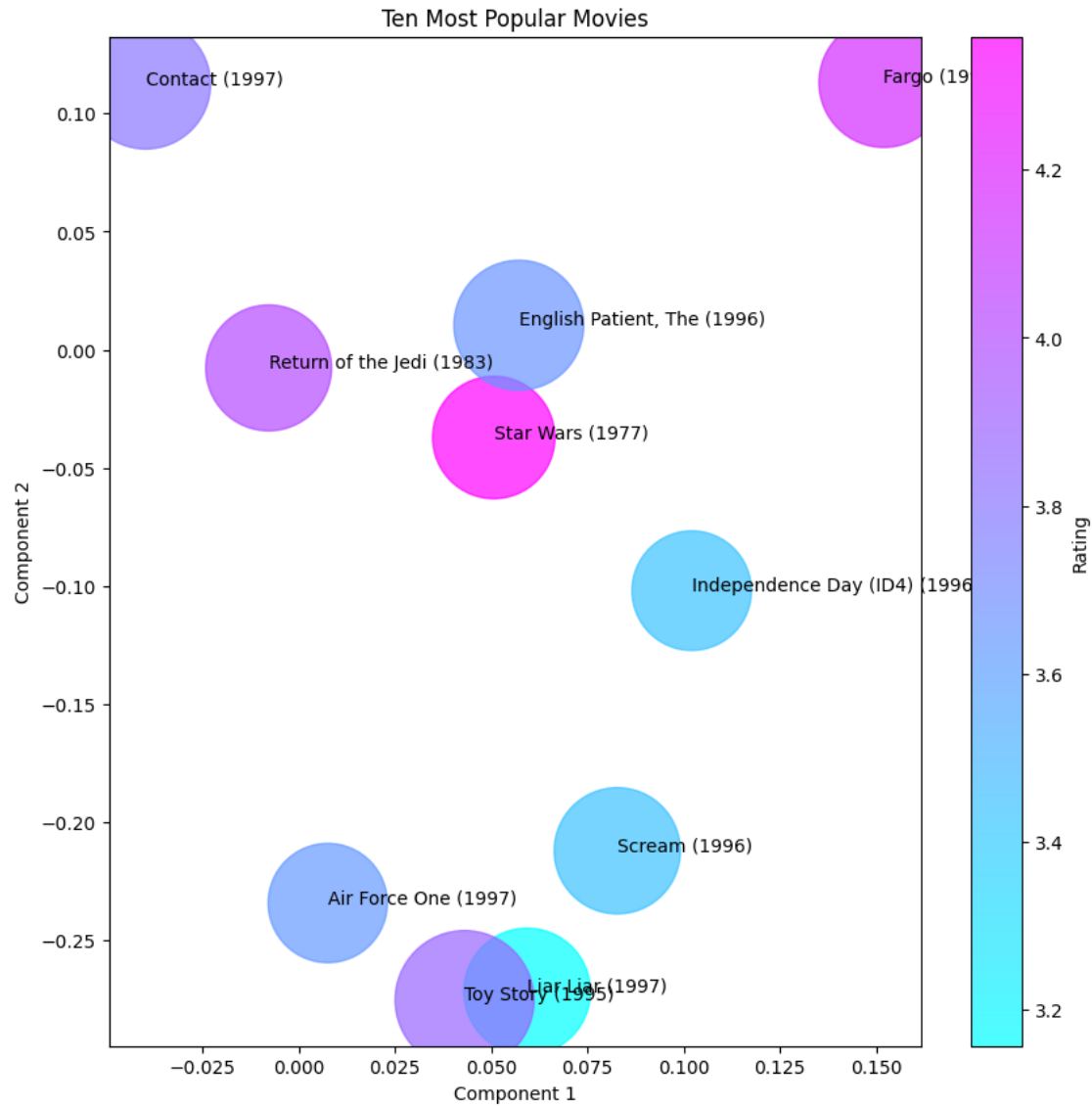


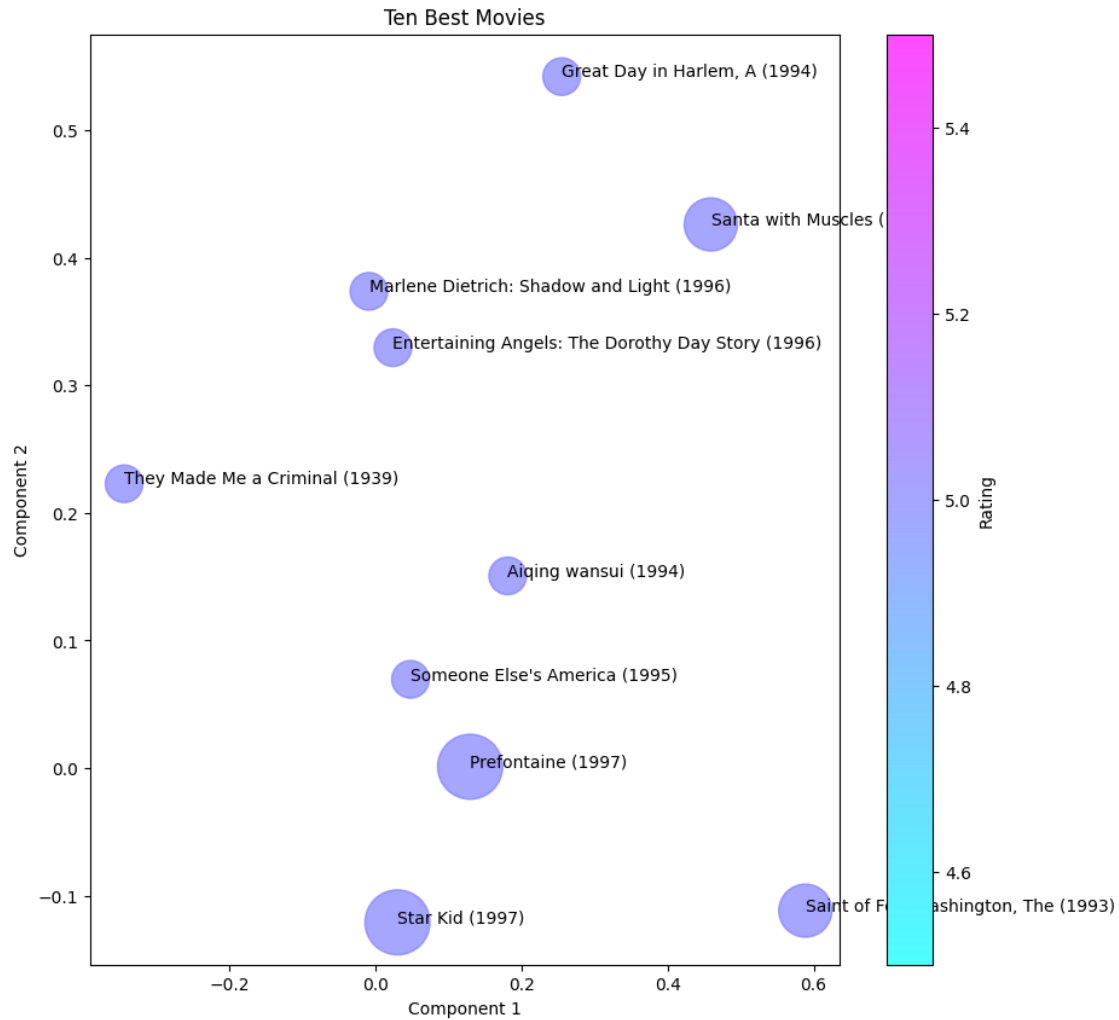


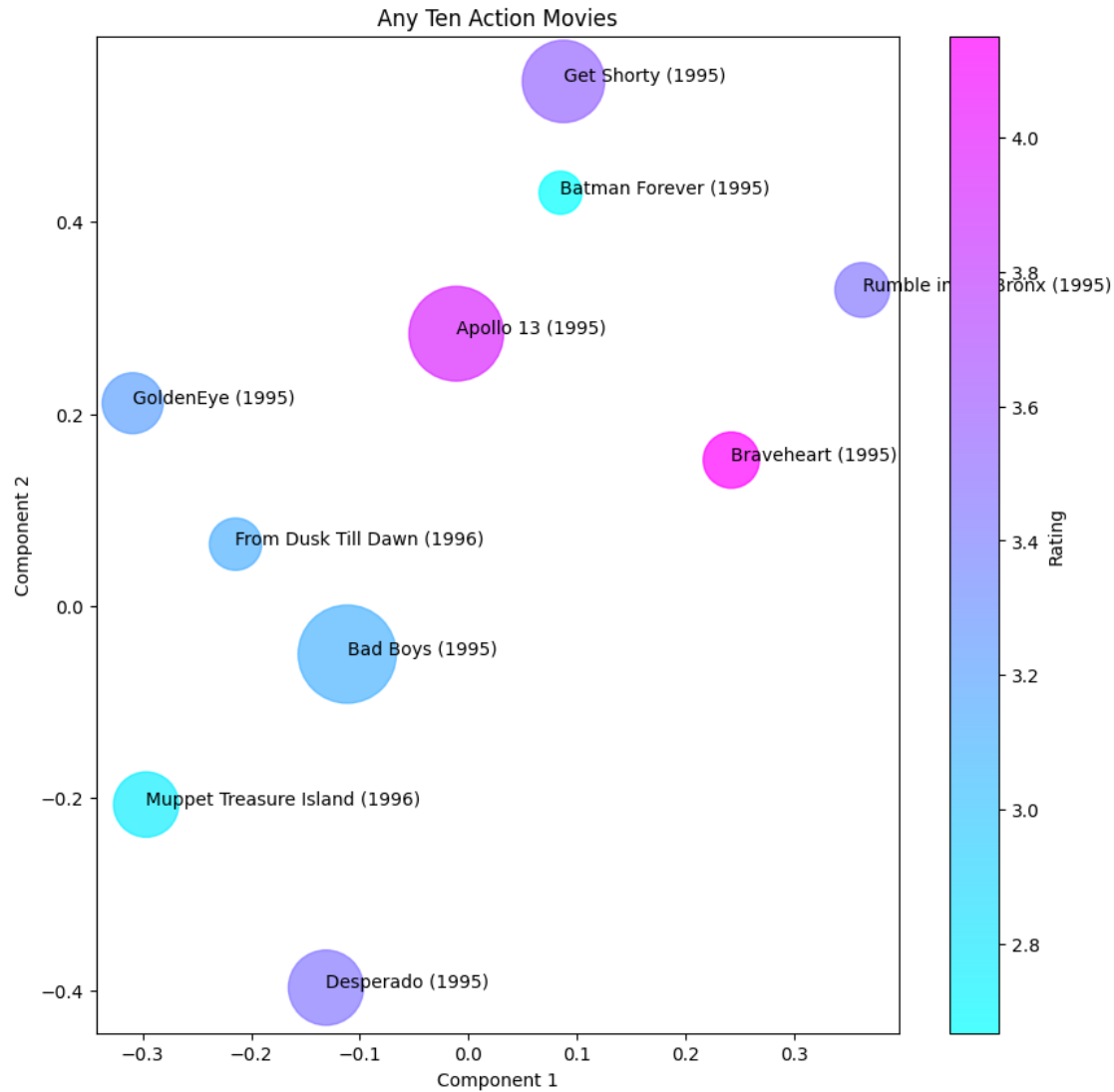
Method 2 Plots

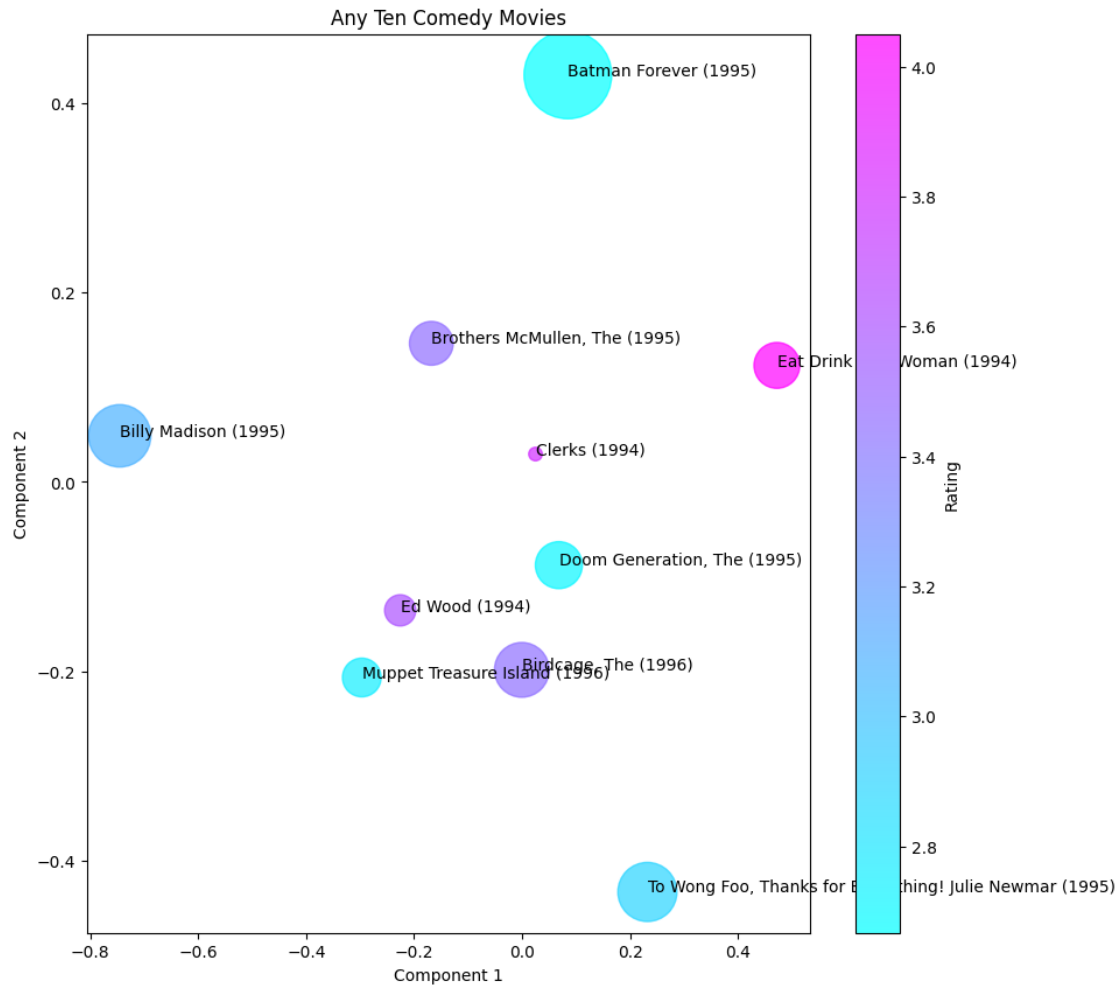
Below are the six visualizations resulting from Method 2.

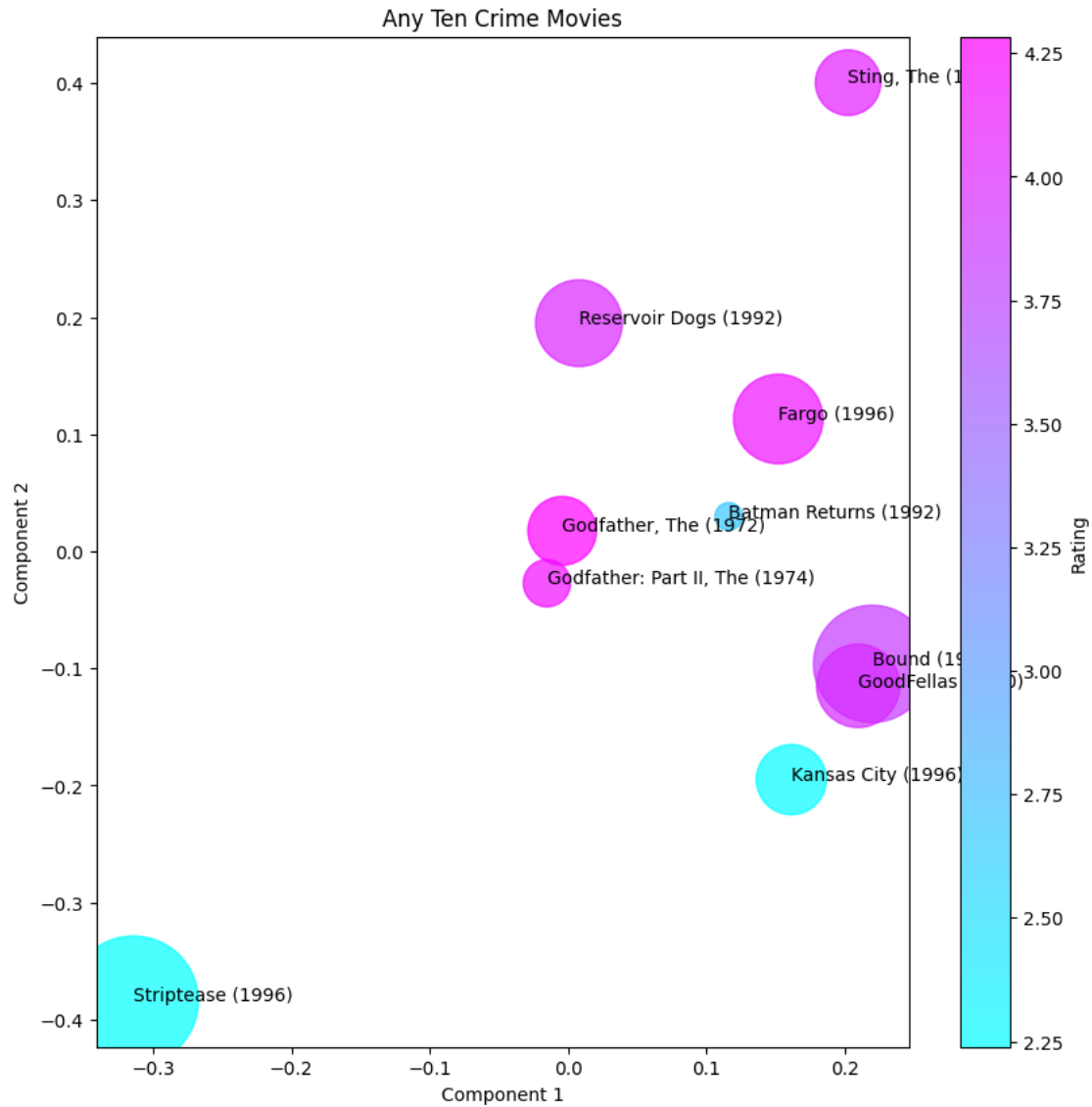






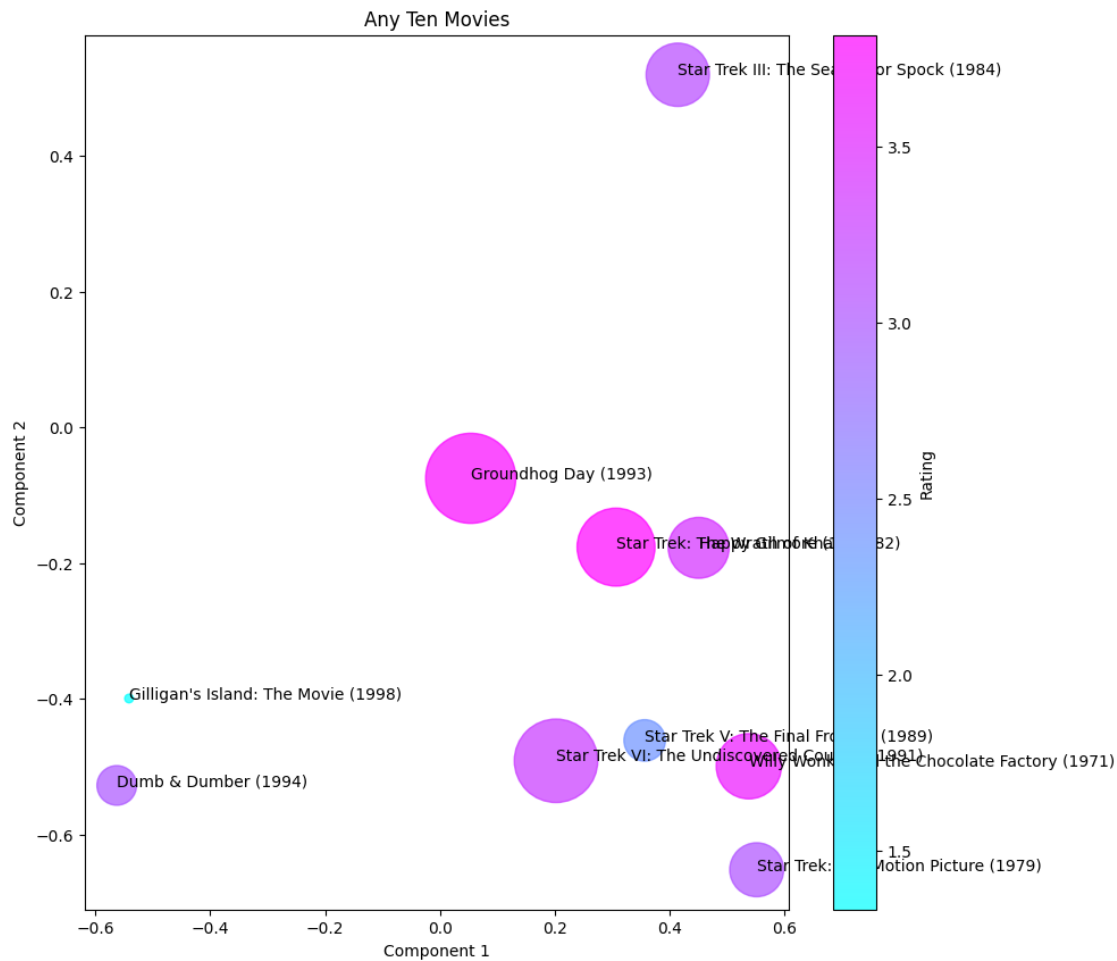


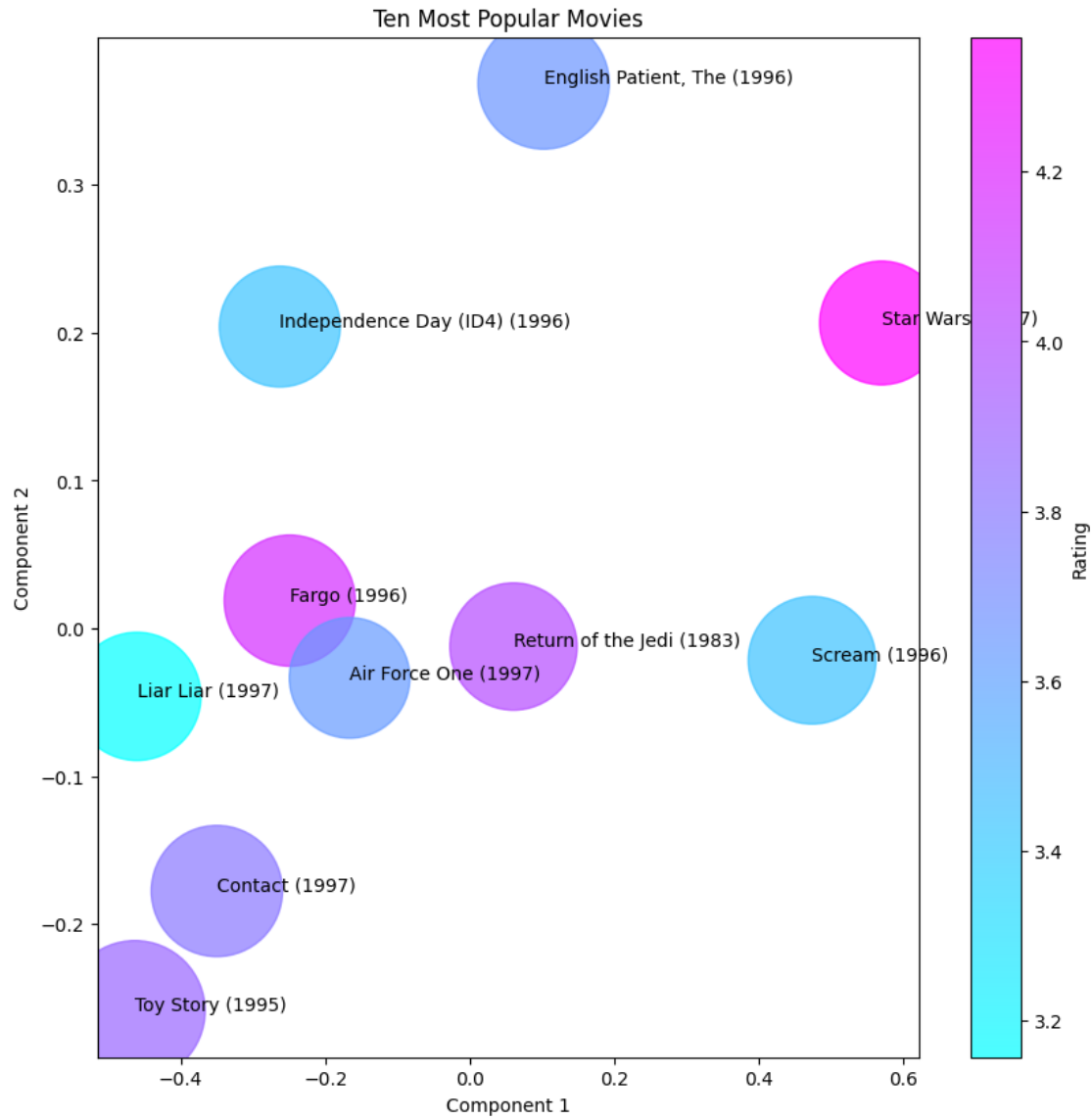


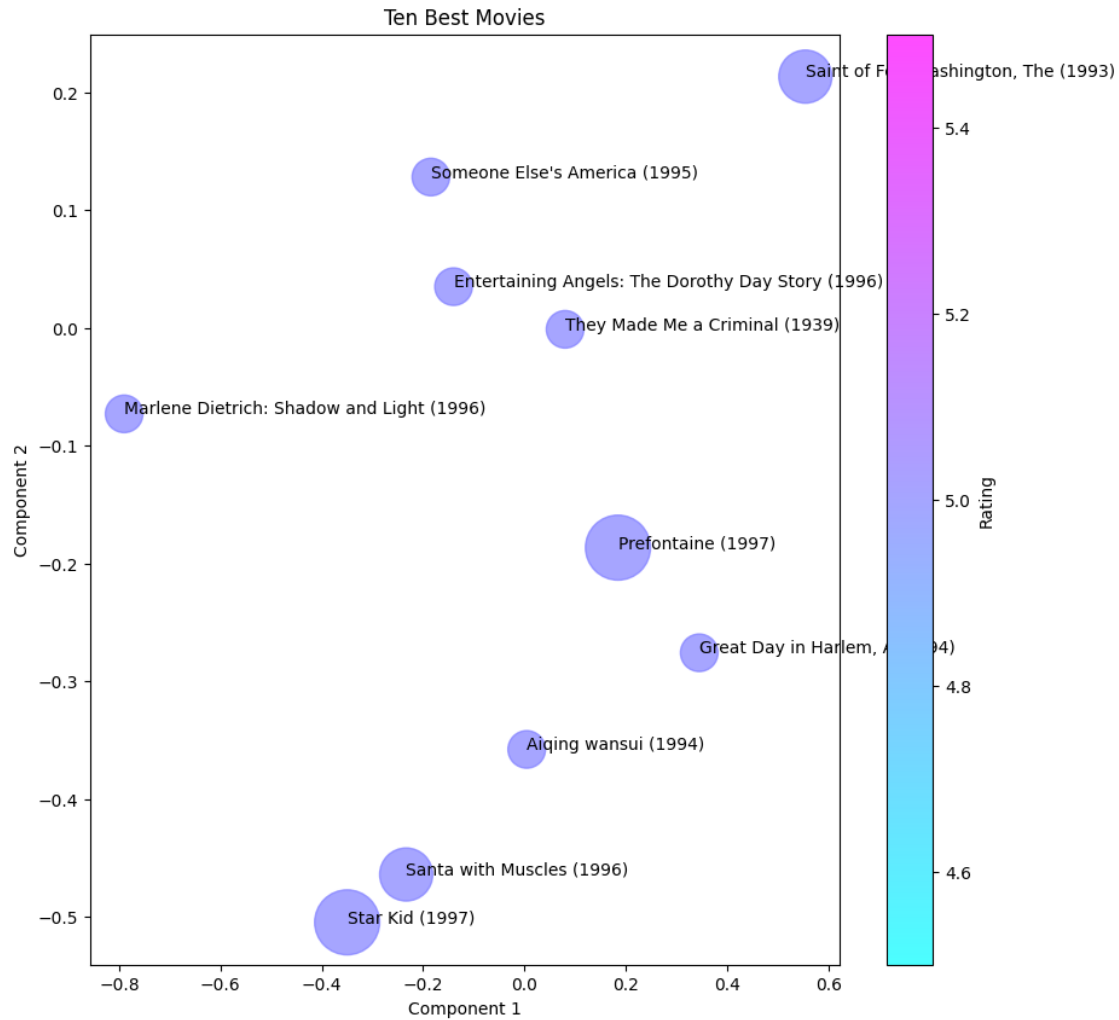


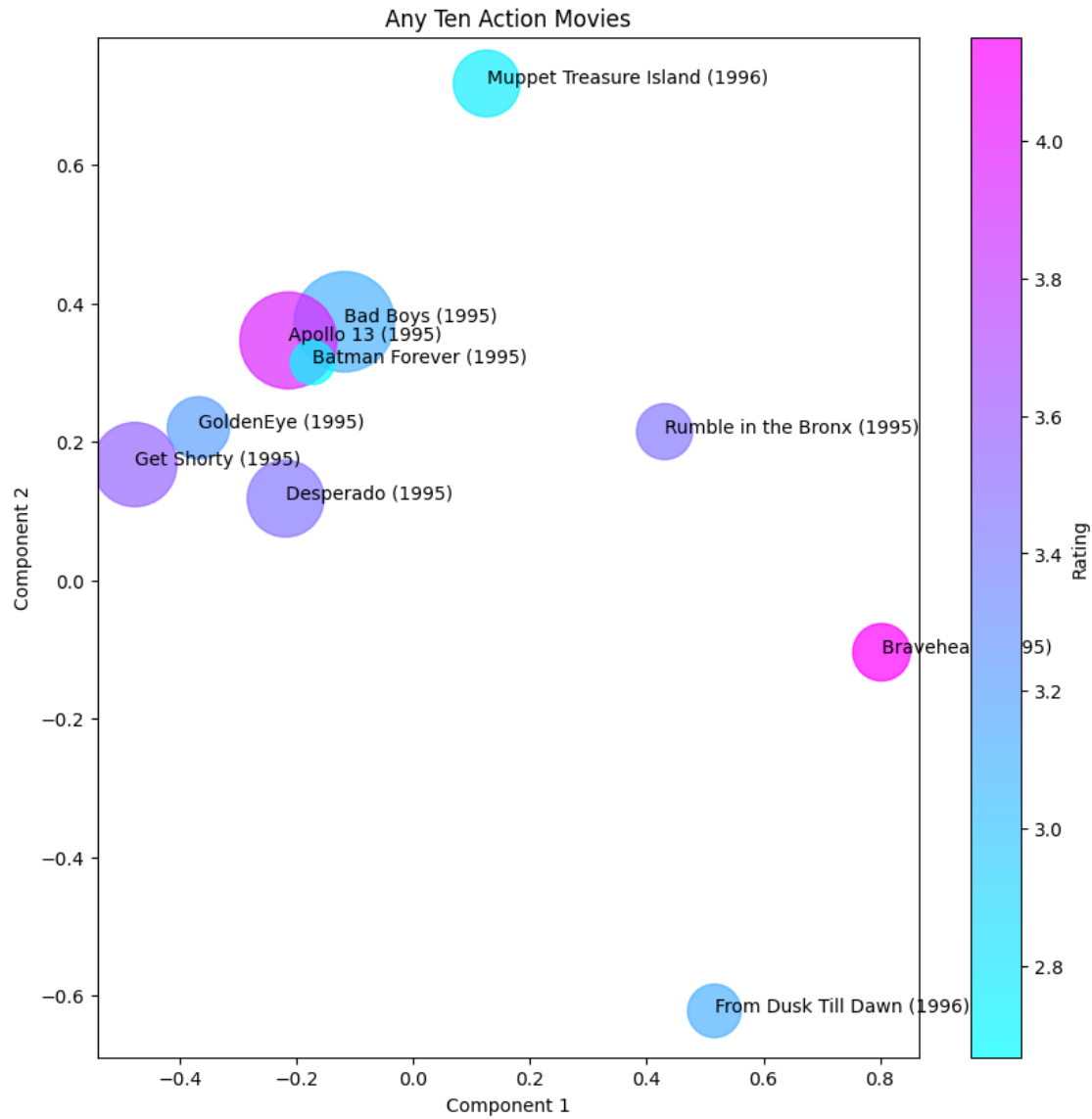
Method 3 Plots

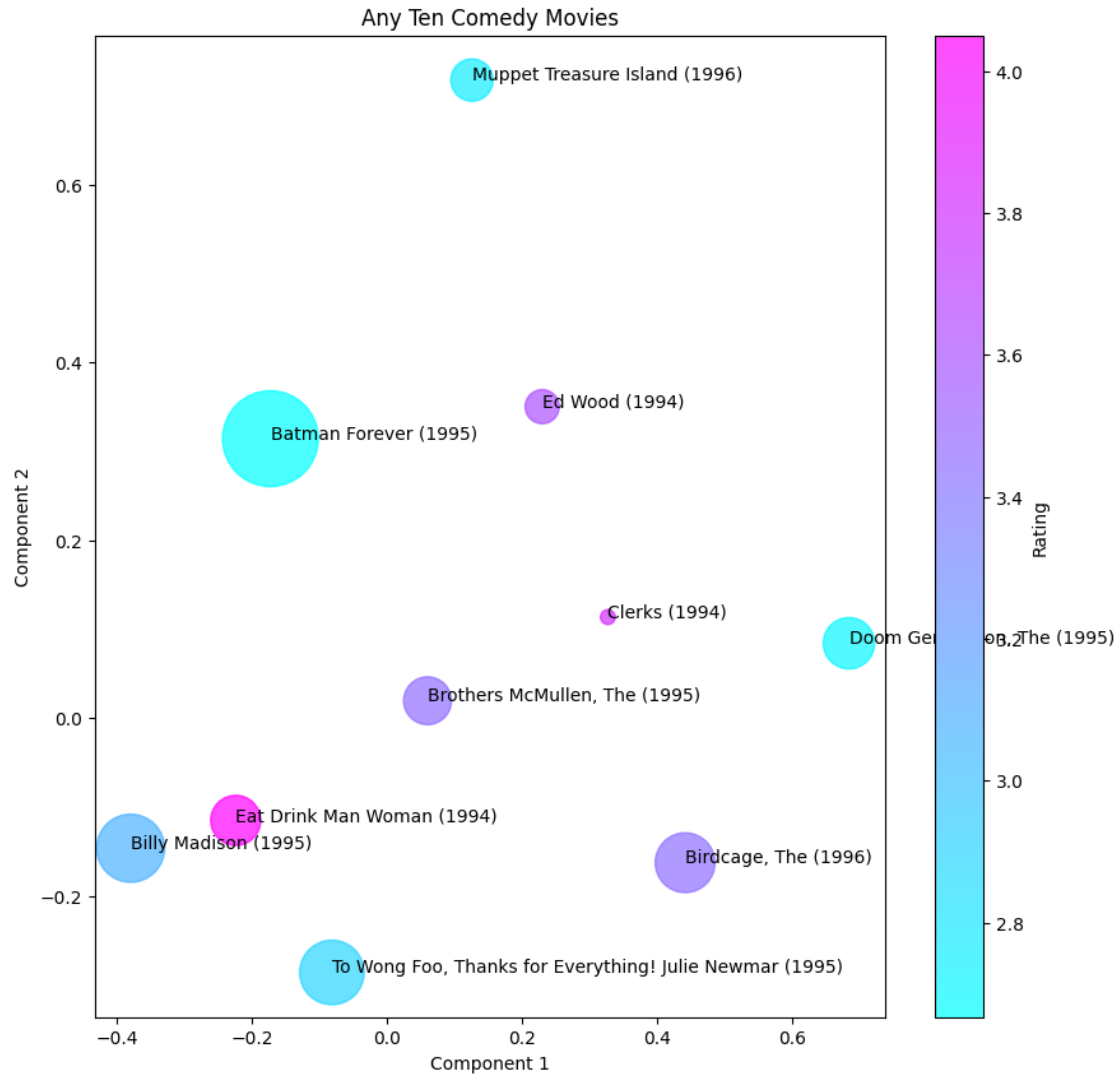
Below are the six visualizations resulting from Method 3.

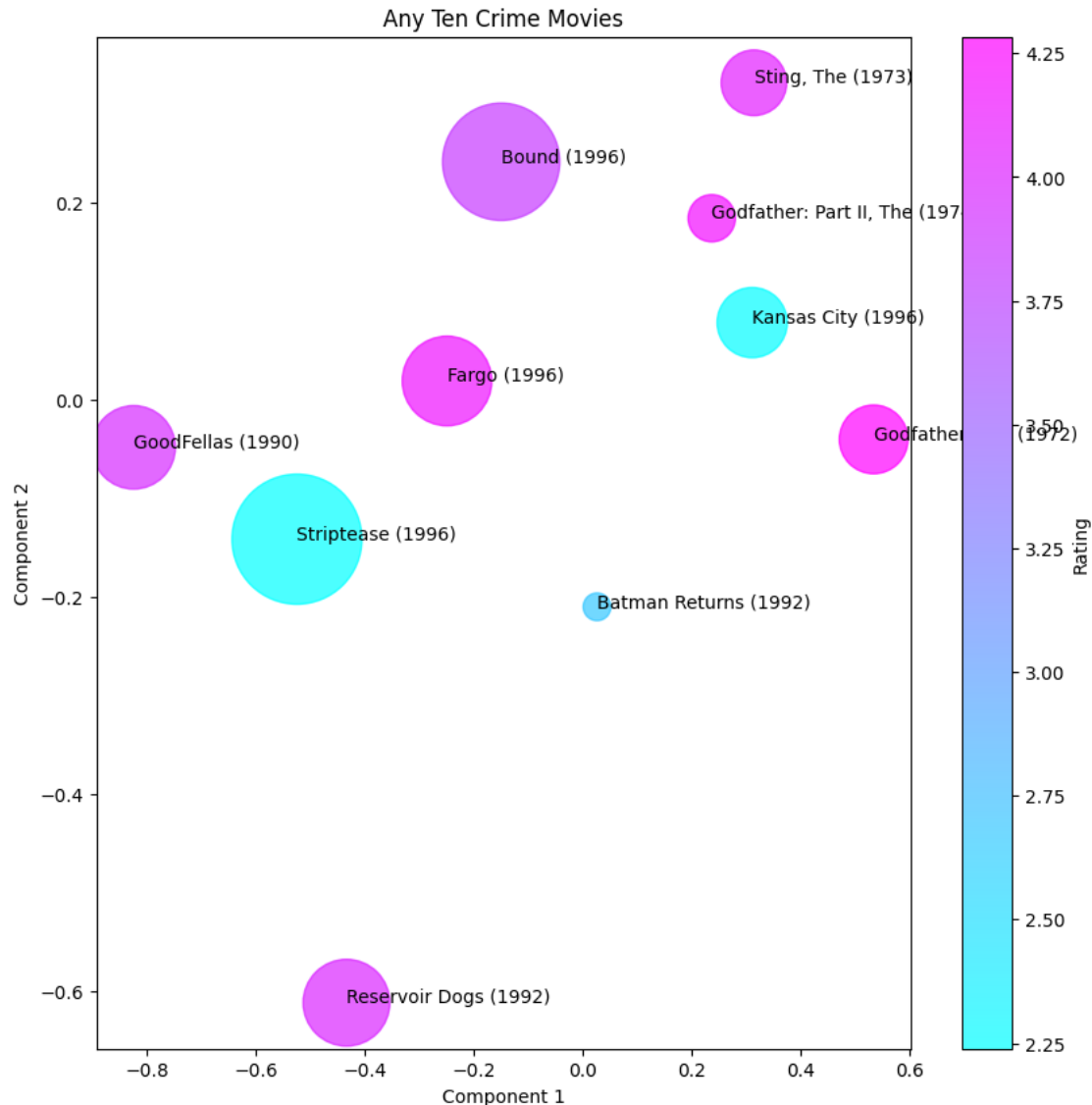












Conclusions from plots

First, we expected to not see any kind of strong trend when it came to the visualization of the best movies. There aren't many reviews included for the best movies so a development of a strong trend would be difficult and not very likely.

For the visualization of the most popular movies, we expected some sort of general trend due to the large amount of ratings. For method 1, we notice a general trend where movies in the top right were the popular movies with relatively lower ratings (around 3 and colored blue in the visualization) and movies with

higher relative ratings were towards the left and lower. A line could be drawn from top right to the lower left to show this trend. The use of colors in the visualization make it easily visible. In method 2 and method 3, the overall presence of movies along this diagonal is followed but the trend with ratings is not. Combining these two, we can suggest that popular movies are more closely related to the values of the upper right and lower left regions of the graph. There again was no observed trend in the distribution of the best movies in methods 2 and 3.

The visualizations of the three genres also provide interesting insights. Many of the action movies were clumped together in Method 3. Other methods had less observed clustering and this may just be due to chance of randomly selecting ten action movies. The other genres had no observed trends when analyzed by Method 3. In method 1, a trend can be seen in all three genres. Movies with lower ratings are farther right on the graph, while movies with higher ratings are more likely to be on the left side of the graph pointing to positioning based on Component 1. Additionally, the crime genre visualization from Method 1 shows clustering by rating in the lower left corner specifically. Compared to the other methods, method 1 seems to show the most obvious trends and clustering. In method 2, the genres, comedy and action, do not seem to show many trends. For crime movies, all movies except one are on the right side of the graph. This could be indicative of a trend or could just be random based on picking of the movies. In method 3, for action movies, some clustering on the left side of the graph can be seen. This could tell us that this is where the majority of action movies are located. Overall, the plots of the movies of a specific genre were not similar for the methods. Specific movies were not in the same general areas and the different methods led to very different graphs.

The methods all offer valuable information and we are pleased with our visualizations. Our special plot is attached below and an in-depth discussion of it is included in our Piazza post.

