

Baselines - Combating Partisan Homogenization in News Recommendation Systems

Contents

1	Introduction	2
2	Problem Statement	2
3	Possible Solutions	2
4	General Experiment Settings	3
4.1	TF-IDF	3
4.2	GLOVE	3
4.3	BERT	3
5	Baseline 1: How Topic Similarity affects Model Performance	4
5.1	TF-IDF	4
5.2	Glove	5
5.3	BERT	6
6	Baseline 2: Online Learning Setting	6
6.1	TF-IDF	7
6.2	Glove	9
6.3	BERT	10
6.4	SUMMARY	11
7	Baseline 3: How easy is it for the Recommendation System to detect a change in topics	12
7.1	TF-IDF	12
7.2	Glove	13
7.3	BERT	14
7.4	Summary	15
8	Baseline 4: Varying Regularization Strength to remove Spurious Correlations	16
8.1	TF-IDF	16
8.1.1	Homogeneous Users	16
8.1.2	Heterogeneous Users	16
8.2	Glove	17
8.2.1	Homogeneous Users	17
8.2.2	Heterogeneous Users	18
8.3	BERT	19
8.3.1	Homogeneous Users	19
8.3.2	Heterogeneous Users	19
8.4	Summary	20
9	Baseline 5: Learning Rate vs Model Performance	21
9.1	TF-IDF	21
9.1.1	Homogeneous Users :	21
9.1.2	Heterogeneous Users :	21
9.2	Glove	22
9.2.1	Homogeneous Users :	22
9.2.2	Heterogeneous Users :	22
9.3	BERT	23
9.3.1	Homogeneous Users :	23
9.3.2	Heterogeneous Users :	23
10	Baseline 6: Model Performance on mixed Set (Cluster 1 + Cluster 2)	24
10.1	TF-IDF	24
10.2	Glove	26
10.3	BERT	28
11	Baseline 7: Learning Rate Variation for Mixed Cluster	30
11.1	TF-IDF	30
11.1.1	Homogeneous	30
11.1.2	Heterogeneous	32
11.2	Glove	34
11.2.1	Homogeneous	34
11.2.2	Heterogeneous	36
11.3	BERT	38
11.3.1	Homogeneous	38
11.3.2	Heterogeneous	40

1 Introduction

2 Problem Statement

3 Possible Solutions

4 General Experiment Settings

github-url : <https://github.com/karthikshivaram24/Combatting-partisan-homogenization>

4.1 TF-IDF

- News Articles Used : **127344**
- Features: **TF-IDF**
- Dimensionality Reduction : **SVD**
- Clustering Algorithm : **K-Means**
- Number of Clusters : **1000**
- Cluster Pair Filtering
 - Minimum Cluster Size : **450**
 - Minimum Partisan Size : **0.5**(used balanced sampling to make label distributions equal)
- Recommendation Model
 - Logistic Regression
 - SGDClassifier with log loss
- Performance Metrics
 - **Macro** : F1, precision , recall
 - **@K** : F1, precision, recall
- User Preferences:
 - **Homogeneous** : Likes articles of the same partisan score across cluster pair (Likes conservative articles in both clusters)
 - **Heterogeneous** : Likes articles of different partisan score across cluster pair (Likes conservative articles in cluster 1 and liberal articles in cluster 2)

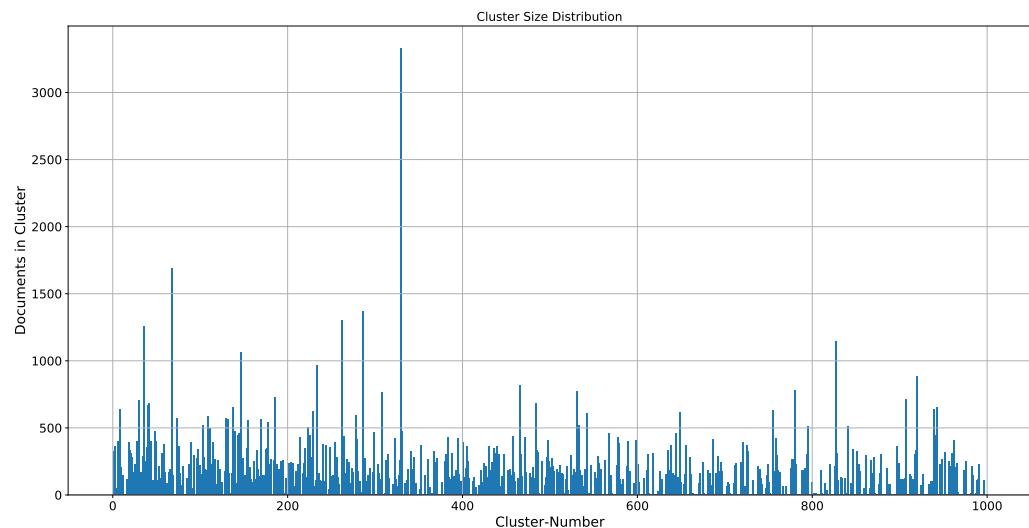


Figure 1: Cluster Size Distribution

4.2 GLOVE

4.3 BERT

5 Baseline 1: How Topic Similarity affects Model Performance

Here we want to measure how well a recommendation system performs on an unseen topic for users with different types of preferences and how this performance varies as similarity between topics increases (seen and unseen topic similarity).

5.1 TF-IDF

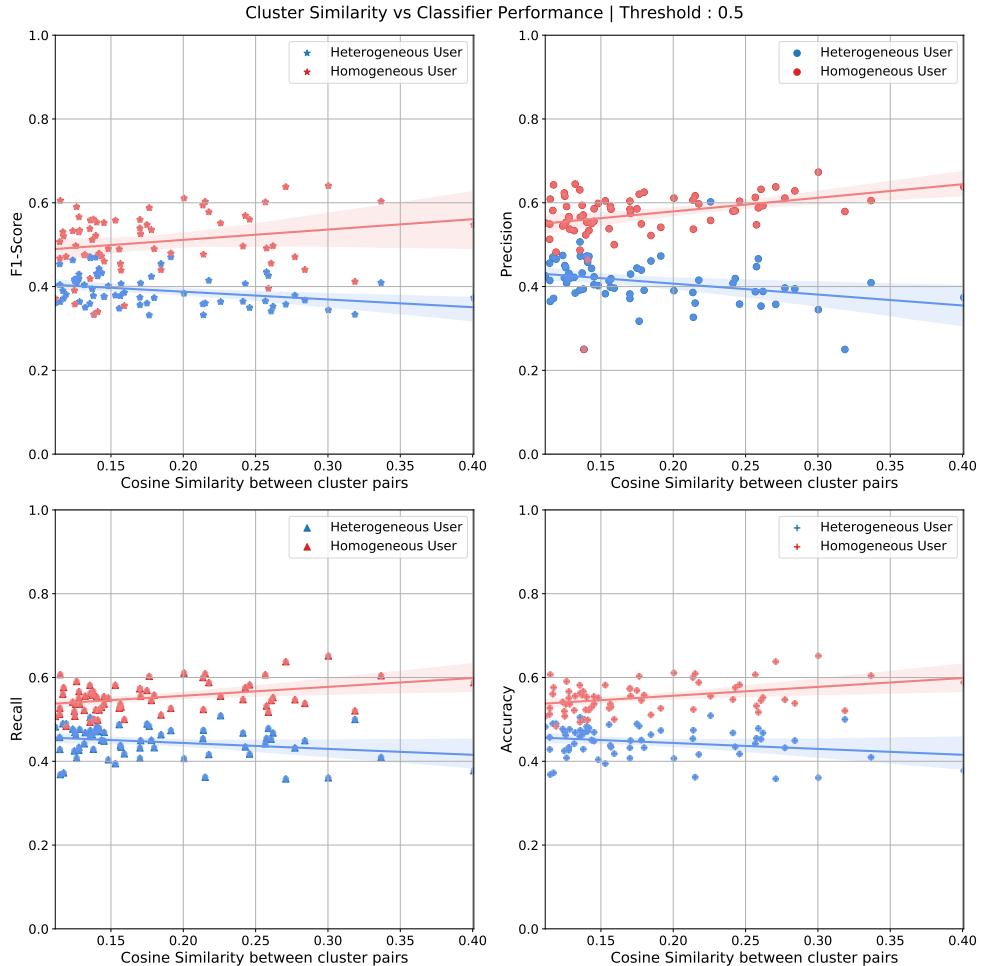


Figure 2: Cosine Similarity between Cluster/Topic Pairs vs F1, Precision, Recall and Accuracy. The models are trained on cluster 1 and metrics are measured on cluster 2 which is used as a test set - TFIDF

Note : The metrics measured for this baseline experiment is based on a test set and not calculated at a particular time step or user-interaction

- We see that for homogeneous user's there is an increase in **Precision** and **Recall** when similarity between the clusters/topics increases.
- For a Heterogeneous User we see the opposite effect with a decrease in both **Precision** and **Recall** as similarity between cluster pairs increases (hence harder for the classifier to distinguish between positives from cluster 1 and negatives from cluster 2 when both are more similar).

5.2 Glove

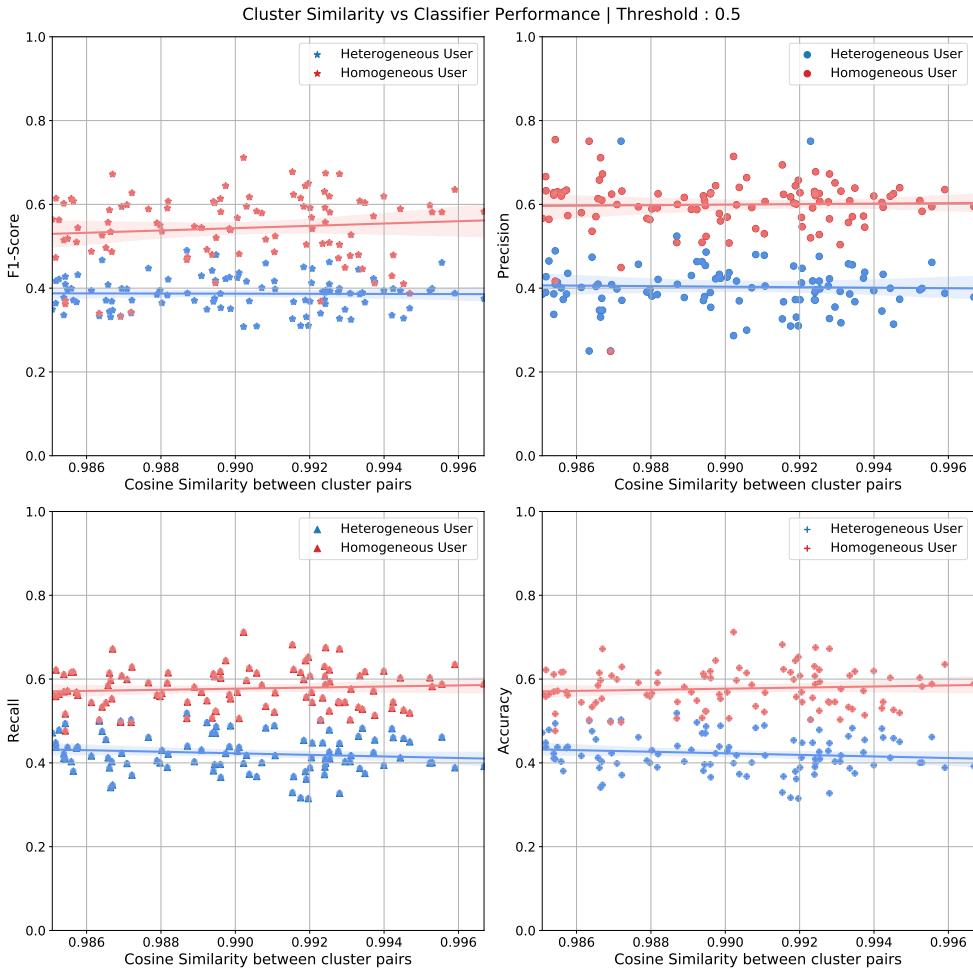


Figure 3: **Cosine Similarity between Cluster/Topic Pairs vs F1, Precision, Recall and Accuracy.** The models are trained on cluster 1 and metrics are measured on cluster 2 which is used as a test set - GLOVE

- We see that for homogeneous user's there is an increase in **Precision** and **Recall** when similarity between the clusters/topics increases (only slight increase occurs here as the cluster pairs seem be more similar here compared to clustering using TFIDF).
- For a Heterogeneous User we see the opposite effect with a decrease in both **Precision** and **Recall** as similarity between cluster pairs increases (but smaller decrease compared to using TFIDF).

5.3 BERT

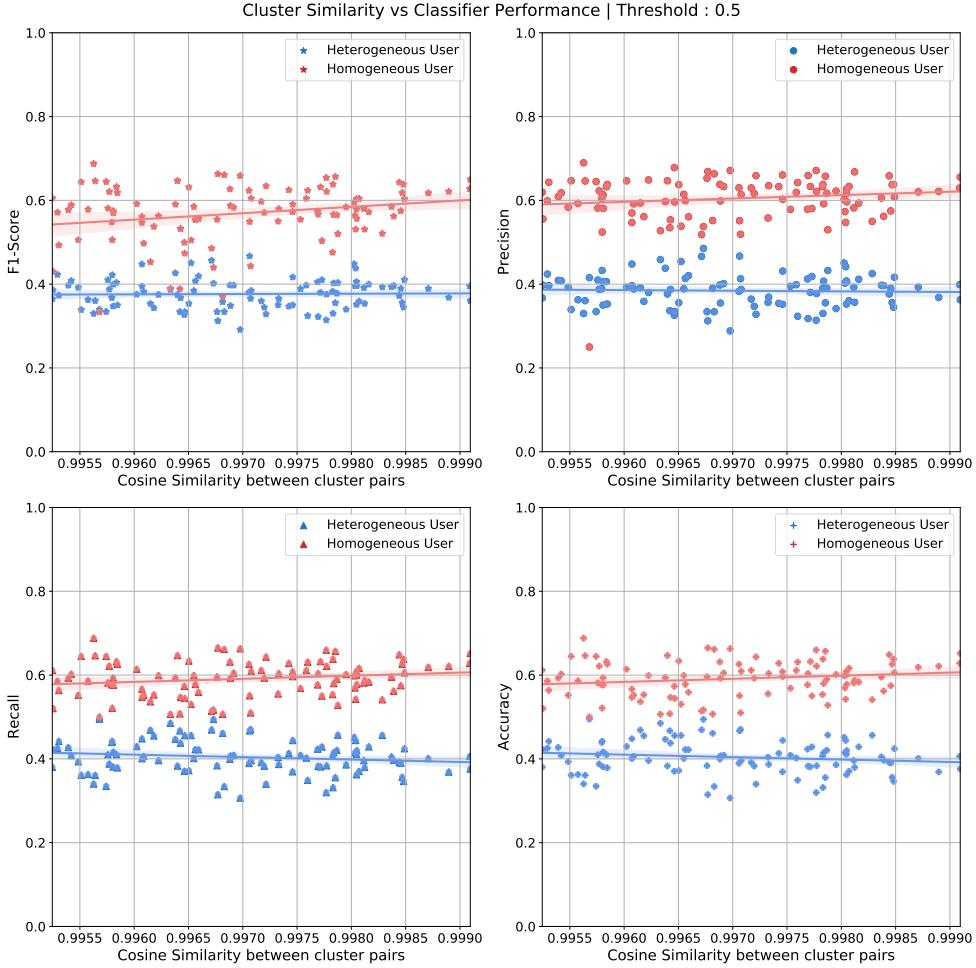


Figure 4: **Cosine Similarity between Cluster/Topic Pairs vs F1, Precision, Recall and Accuracy.** The models are trained on cluster 1 and metrics are measured on cluster 2 which is used as a test set - BERT

- Similar Performance like GLOVE but the clustering of pairs is even better now as we choose extremely similar cluster pairs measured according to cosine similarity.

6 Baseline 2: Online Learning Setting

To emulate a real-world scenario , we want to simulate a recommendation system interacting with a user over a set of unseen articles , slowly updating itself to learn the user's preferences over time, $N=200$ is used here to calculate the metrics, where N represents the number of interactions between the user and the content recommender. Also to note, there are at least N relevant articles in the candidate pool.

For this baseline experiment we use recommender specific metrics such as **Precision@K** and **Recall@K**.

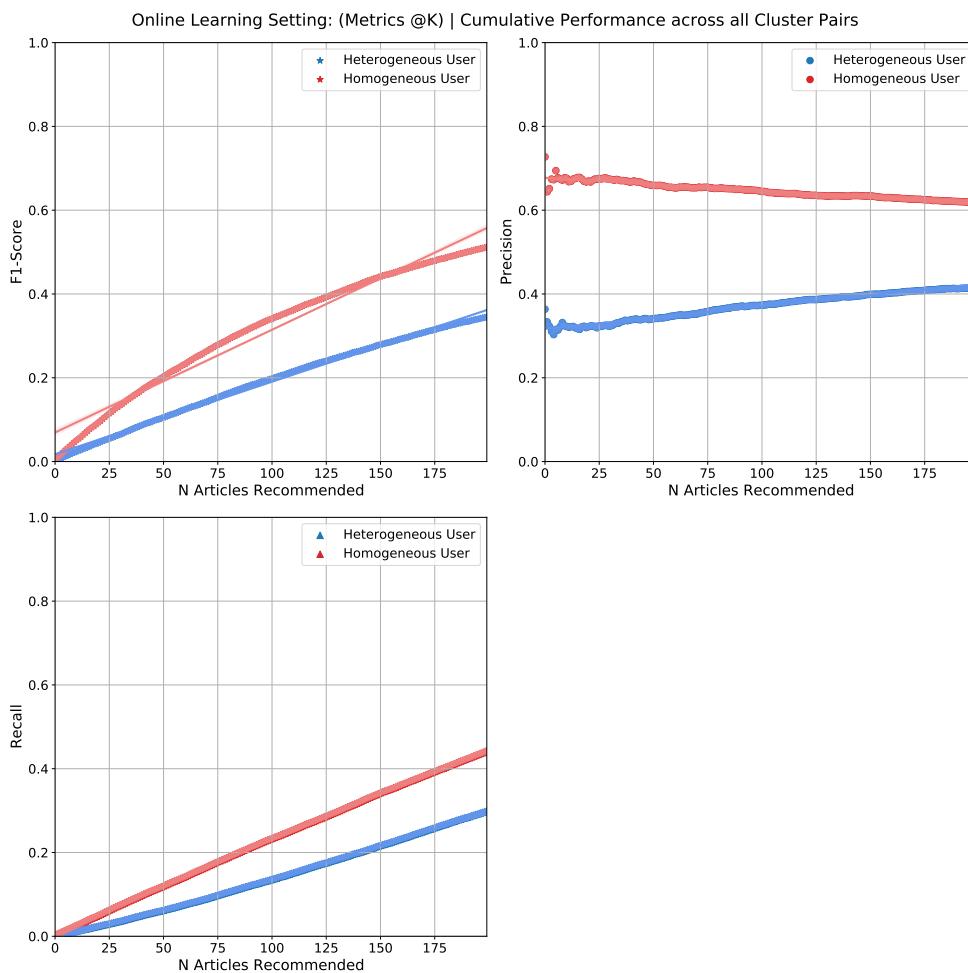
- $Precision@K = \frac{\#recommended\ items\ the\ user\ likes\ at\ K_{th}\ Interaction}{total\ \#\ of\ items\ recommended\ at\ K_{th}\ Interaction}$
- $Recall@K = \frac{\#recommended\ items\ the\ user\ likes\ at\ K_{th}\ Interaction}{total\ \#\ of\ relevant\ items\ in\ the\ candidate\ pool}$

• Eg:

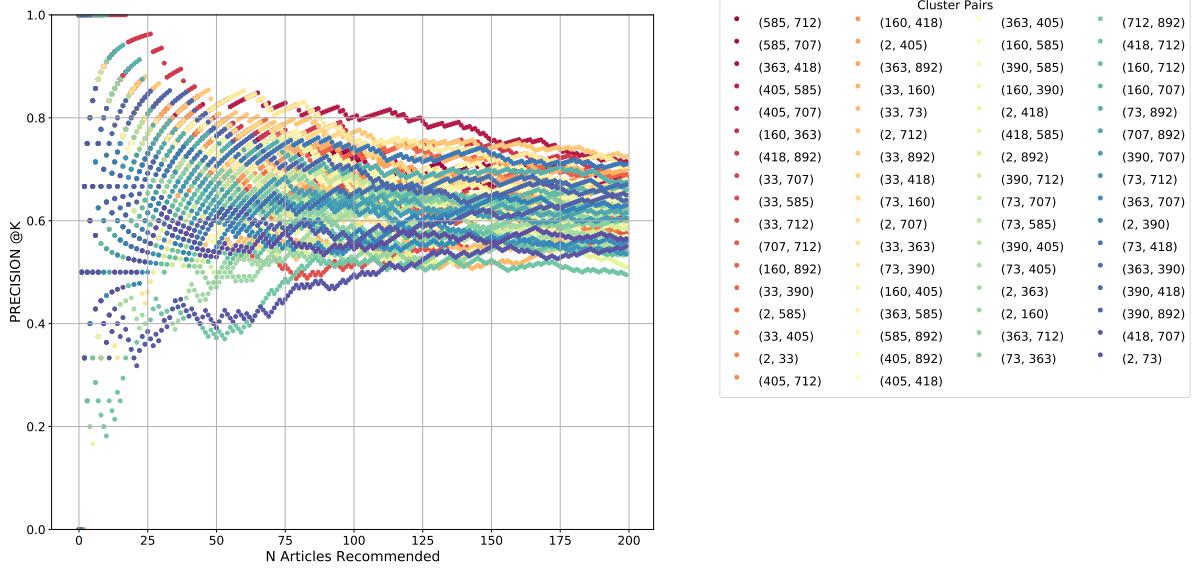
- Candidate Pool Size = 15
- Number of relevant(user-likes) items in the candidate pool = 10
- $N = 1$
 - * Shown Articles = [0]
 - * Recall = $0/10 = 0$
 - * Precision = $0/1 = 0$
- $N = 2$
 - * Shown Articles = [0,1]
 - * Recall = $1/10 = 0.1$

- * Precision = $1/2 = 0.5$
- N = 3
 - * Shown Articles = [0,1,1]
 - * Recall = $2/10 = 0.2$
 - * Precision = $2/3 = 0.666$
- N = 4
 - * Shown Articles = [0,1,1,1]
 - * Recall = $3/10 = 0.3$
 - * Precision = $3/4 = 0.75$

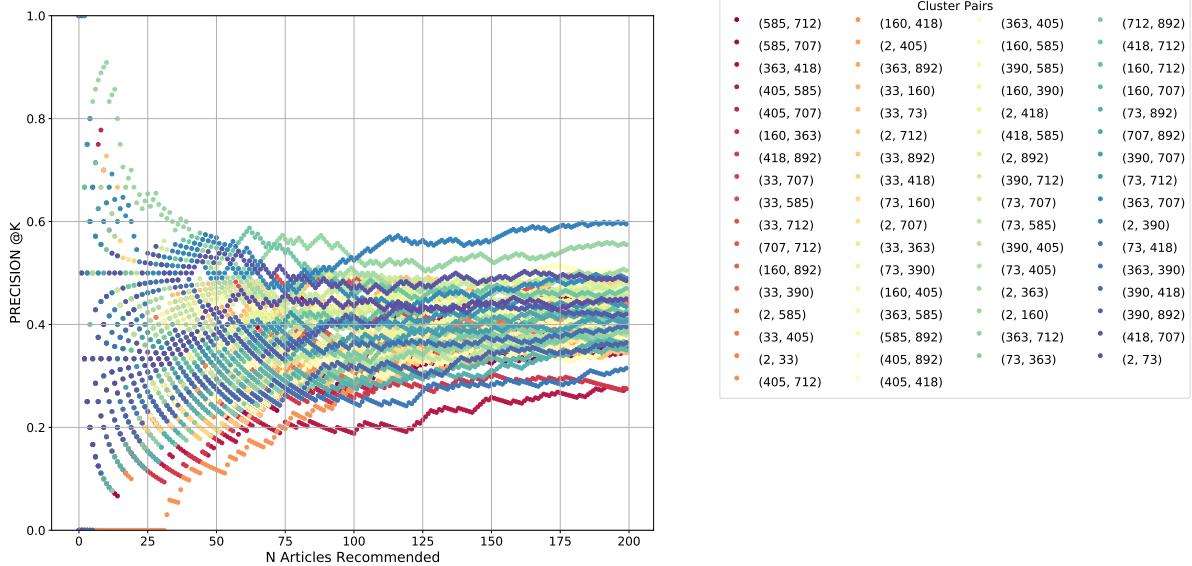
6.1 TF-IDF



PRECISION | Homogeneous



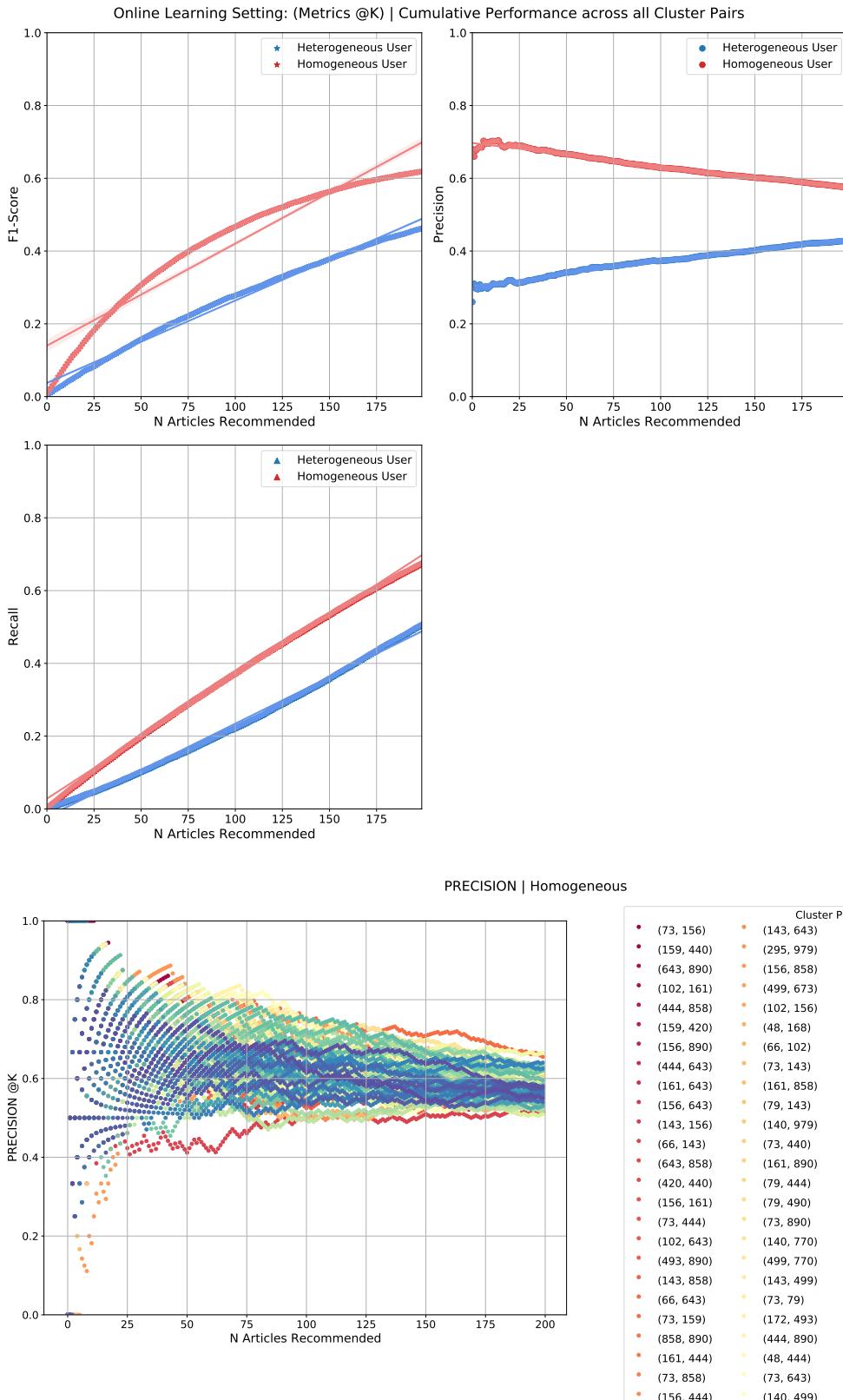
PRECISION | Heterogeneous



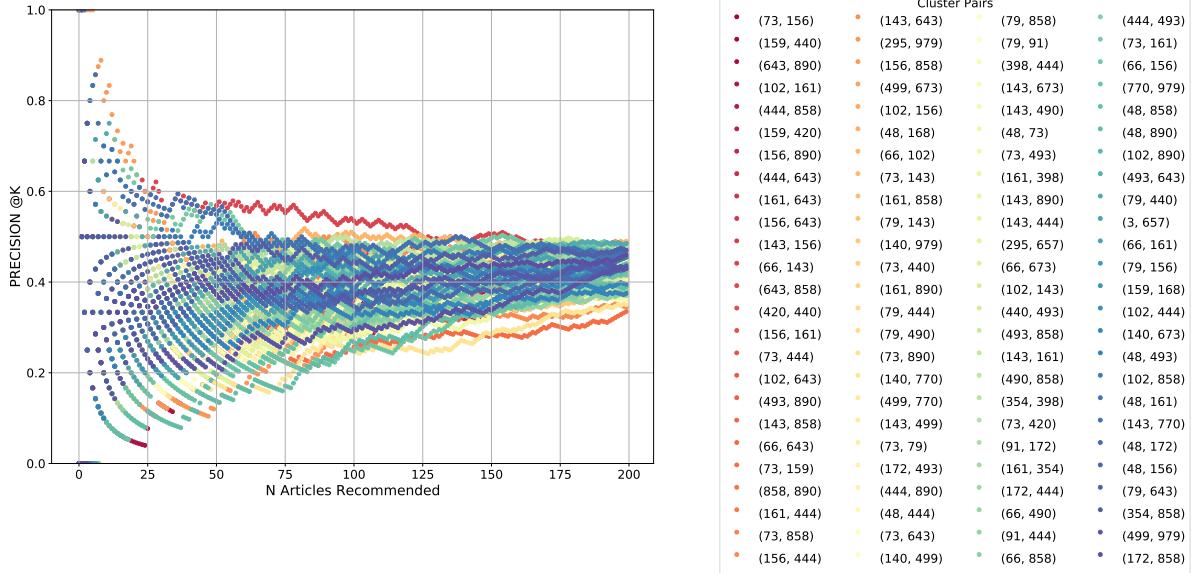
The above figure shows the mean results across all cluster pairs.

- For the Homogeneous User we see that **Precision@K** decreases as the number of interactions increase, this could be due to the fact that the most probable items the user likes are already recommended in the initial user-interactions (as we sort by predicted probability) and as interactions increase the classifier is only left with items with low confidence to recommend hence increasing mistakes. Also from the precision graph with all cluster pairs shown we see that a few cluster pairs do increase in precision.
- For the Heterogeneous User we see that there is an initial dip in **Precision@K** and then it increases as the content recommender learns the different preferences of the heterogeneous user.
- Overall we see that the precision for Homogeneous users is higher than heterogeneous users

6.2 Glove

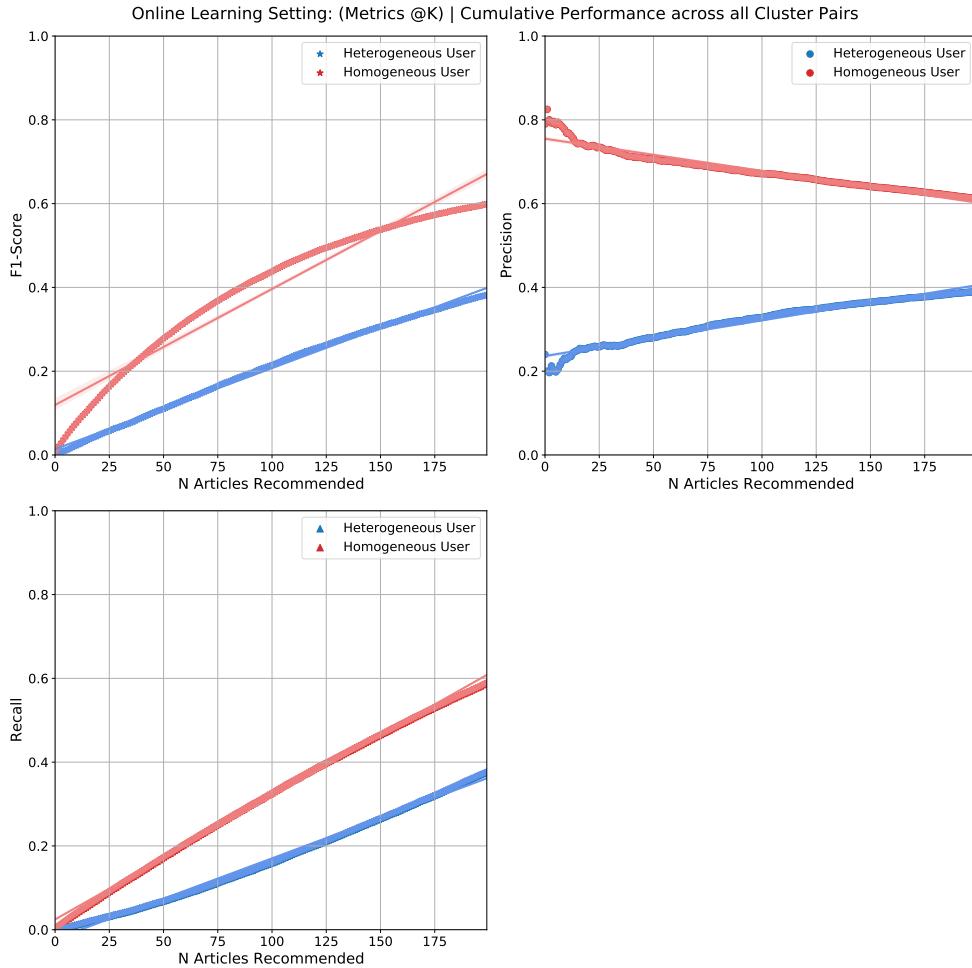


PRECISION | Heterogeneous

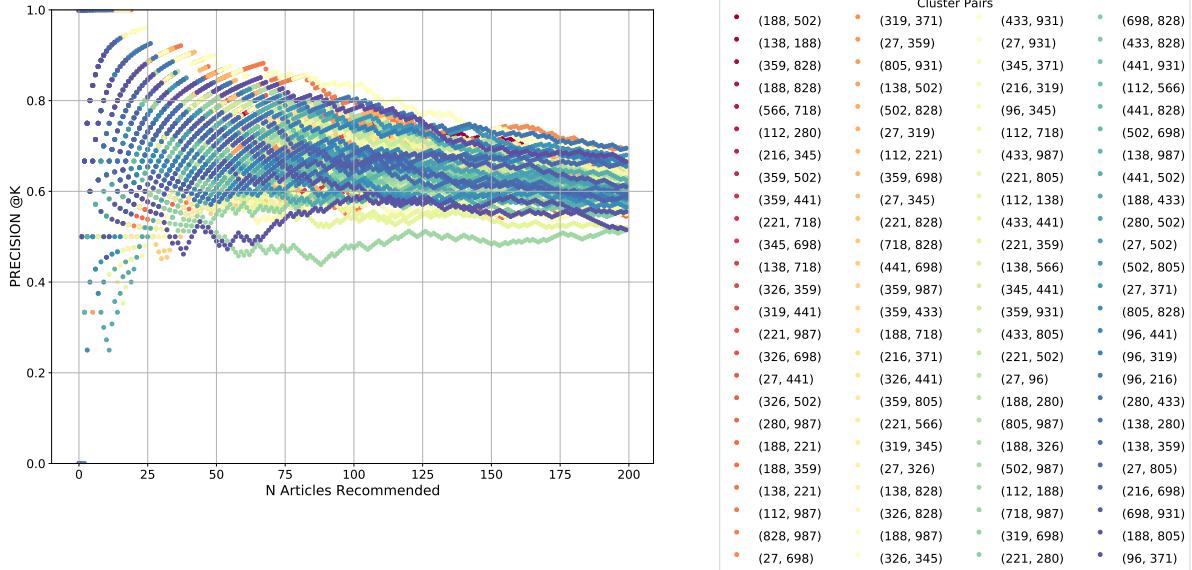


- For Homogeneous Users here we see that their is a larger decrease in Precision compared to using TFIDF representations.
- For the Heterogeneous User we see a larger increase in Precision compared to using TFIDF representations

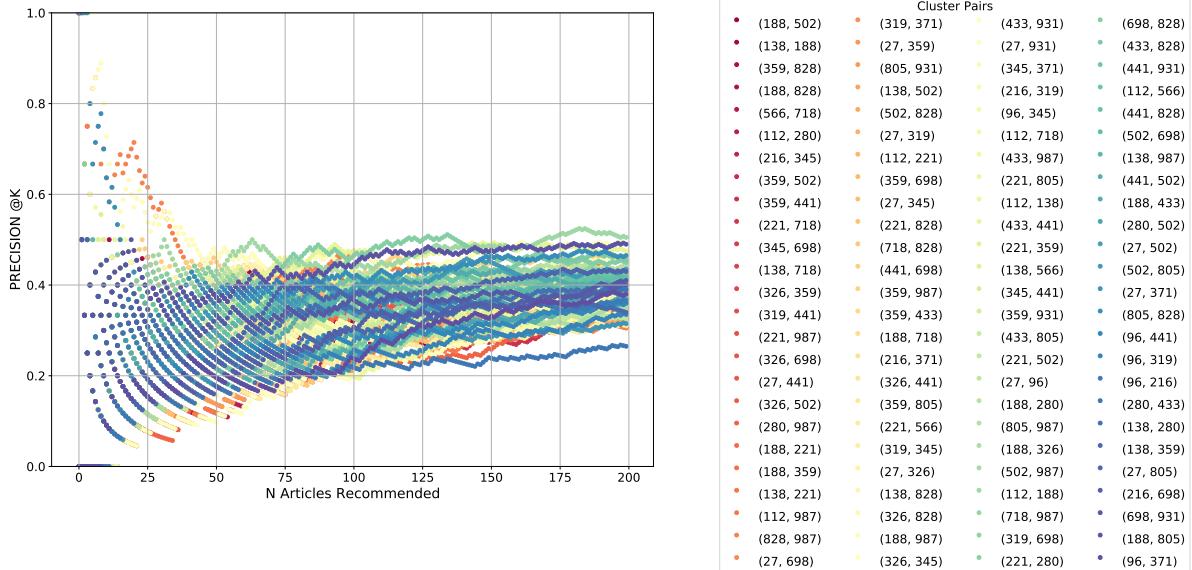
6.3 BERT



PRECISION | Homogeneous



PRECISION | Heterogeneous



- For Homogeneous Users here we see that their is a smaller decrease in Precision compared to using GLOVE representations.
- For the Heterogeneous User we see a smaller increase in Precision compared to using GLOVE representations

6.4 SUMMARY

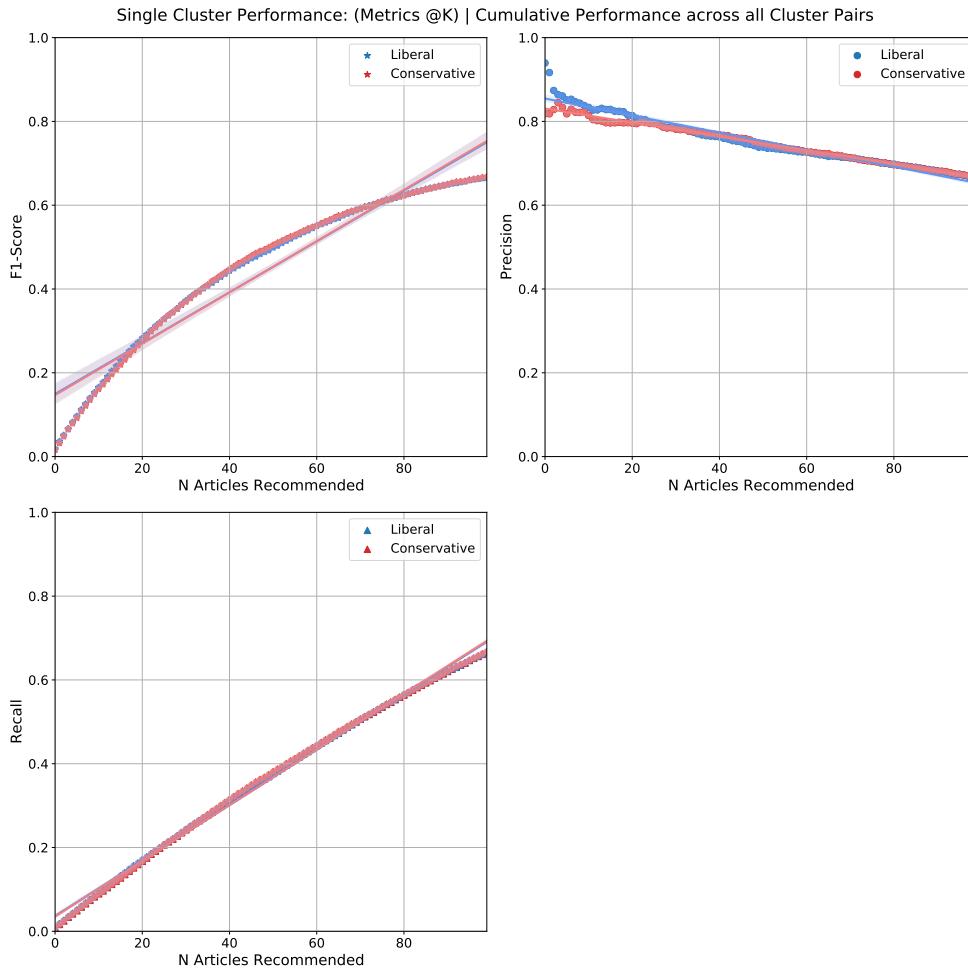
- Using static or non-contextual embeddings helps the system's performance for Heterogeneous Users (according to Precision@K)
- Using dynamic or contextual embeddings helps the system's performance for Homogeneous Users.

7 Baseline 3: How easy is it for the Recommendation System to detect a change in topics

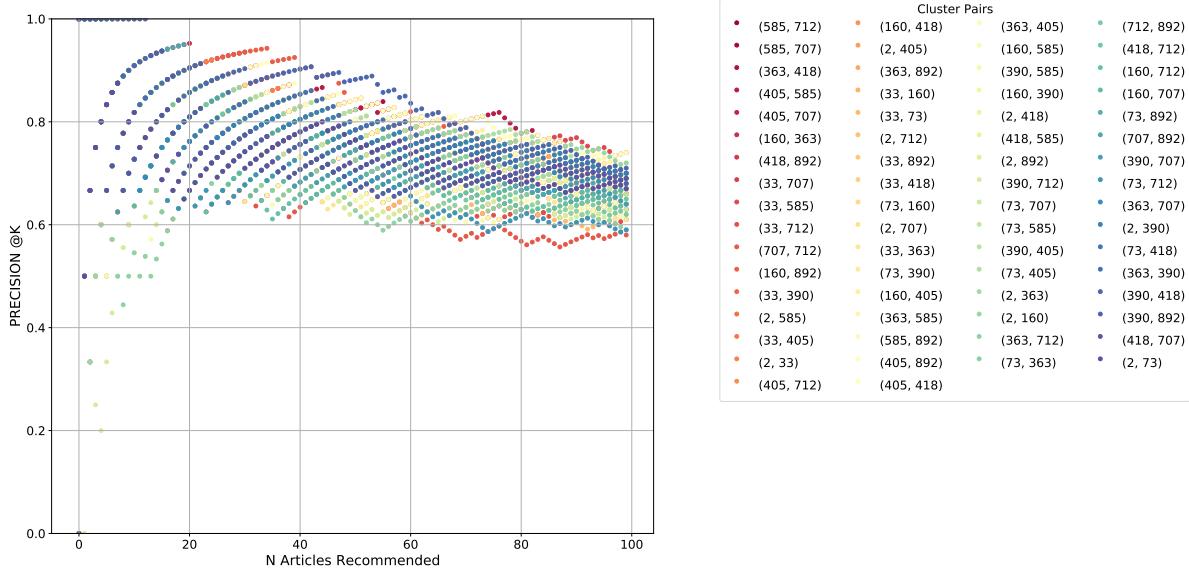
We want to know how the recommendation system performs in detecting a change in topics , so we measure the performance of the system using a single cluster and compare it against our online setting performance (shown in the above baseline).

7.1 TF-IDF

- Similar trend in Precision@K occurs here (compared to baseline 2 - Homogeneous User).
- When we compare these scores to graphs in Baseline 2 we definitely see that the recommendation system tends to have an easier time when only one topic is concerned compared to 2 different topics, so the model is having a hard time in identifying a change in topic.

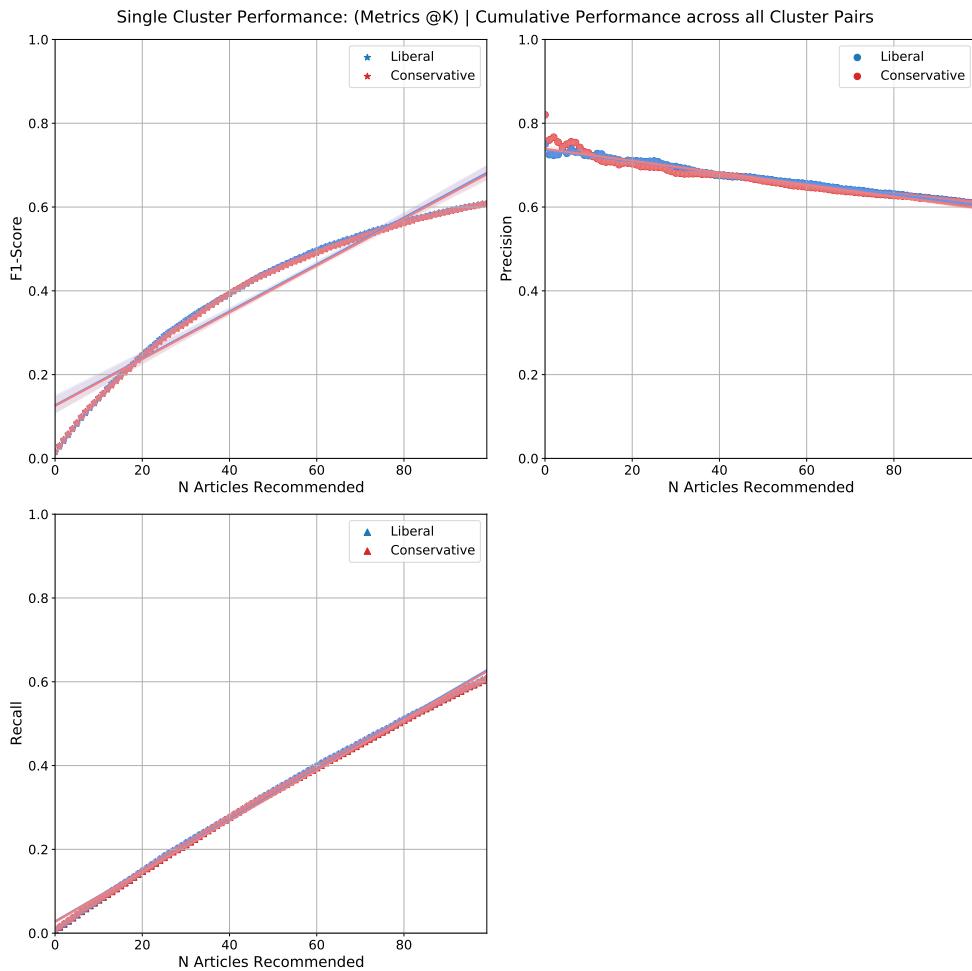


PRECISION | Single Cluster Performance

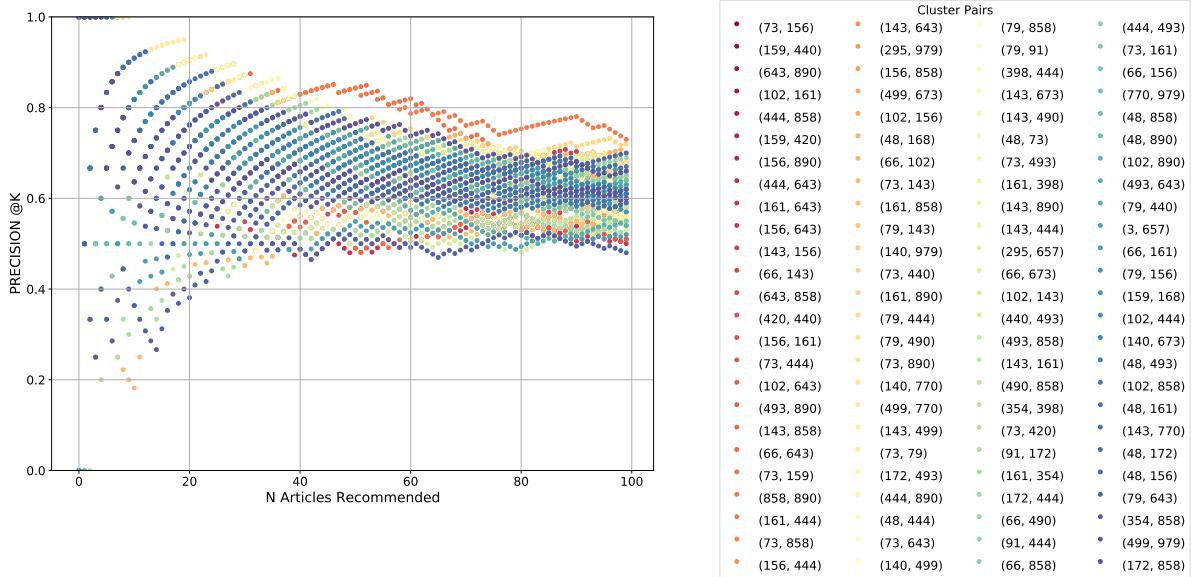


7.2 Glove

- Similar trend in Precision@K occurs here (compared to baseline 2 - Homogeneous User).
- Lower Precision compared to TFIDF Representations

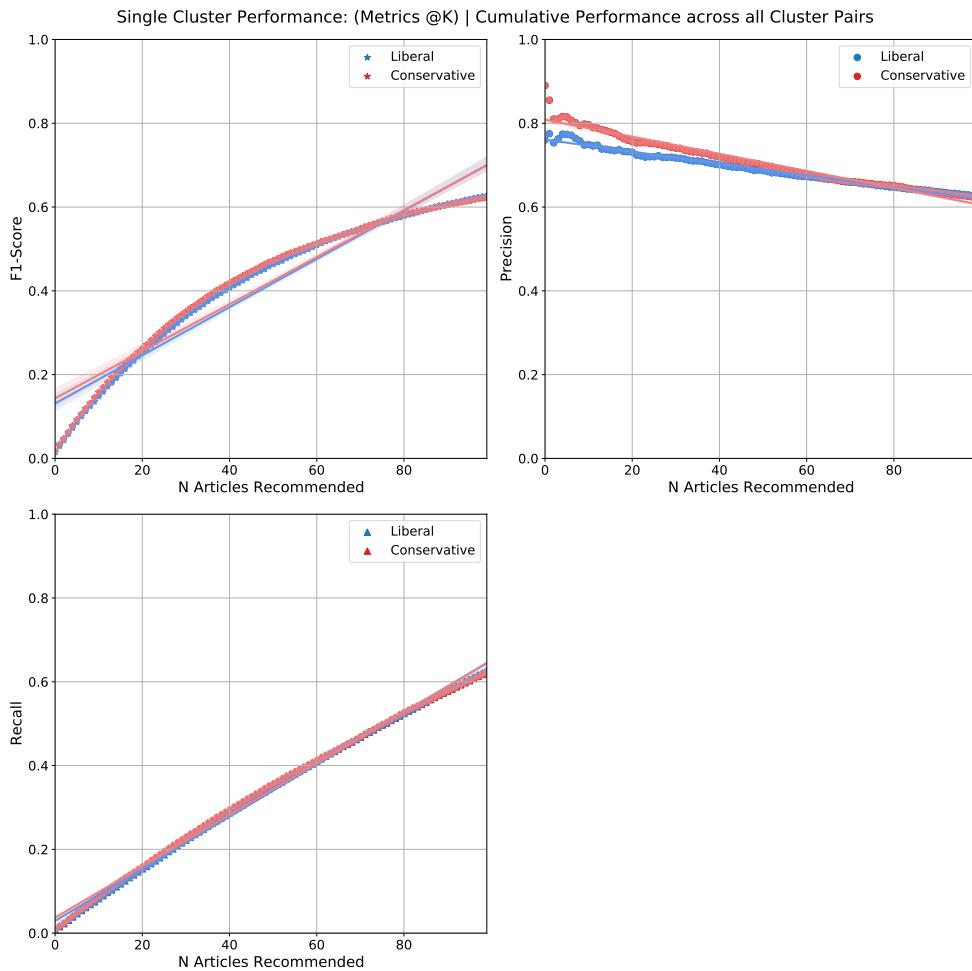


PRECISION | Single Cluster Performance

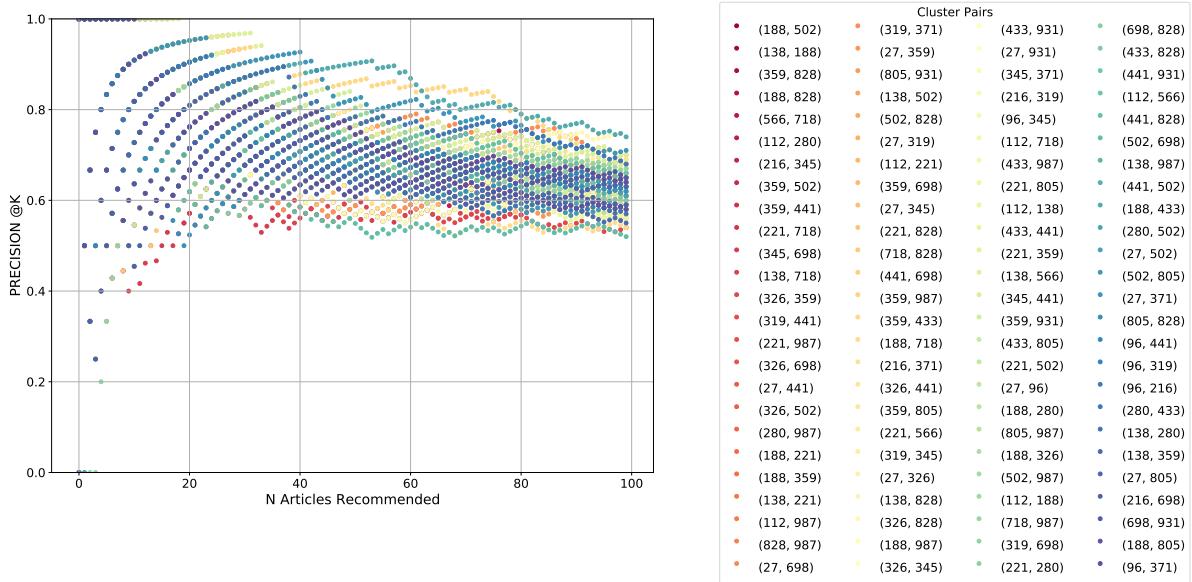


7.3 BERT

- Similar trend in Precision@K occurs here (compared to baseline 2 - Homogeneous User).
- Lower Precision trend compared to TFIDF representations but better than Glove



PRECISION | Single Cluster Performance



7.4 Summary

- Better performance when only one topic/cluster is involved compared to a cluster pair setting
- TF-IDF has the best performance here followed by BERT and then GLOVE

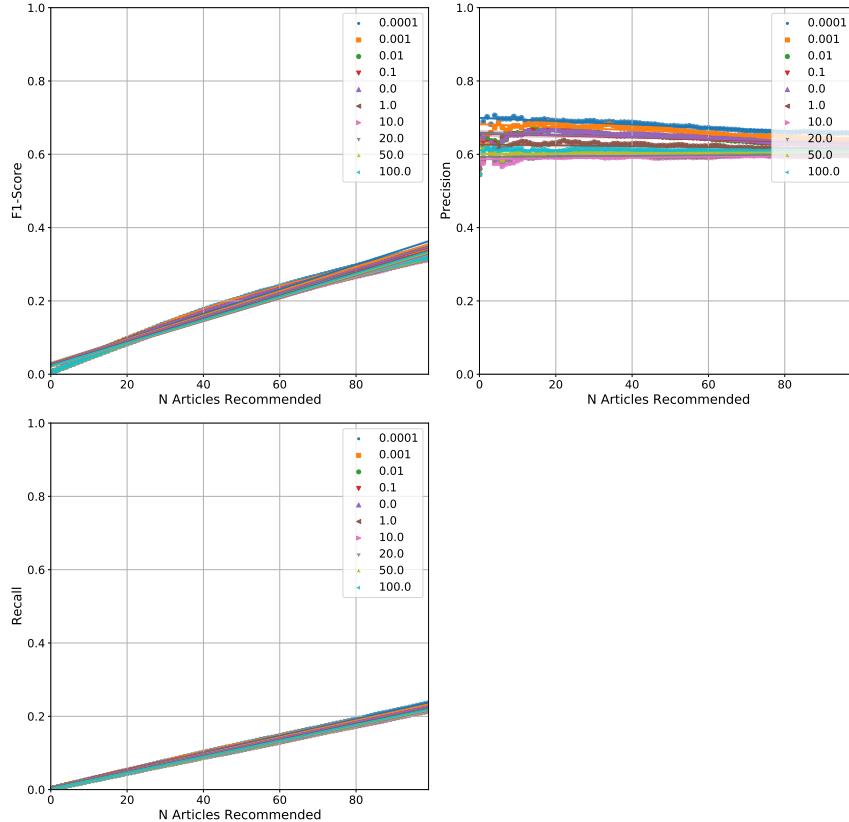
8 Baseline 4: Varying Regularization Strength to remove Spurious Correlations

8.1 TF-IDF

8.1.1 Homogeneous Users

For Homogeneous Users we see that a high regularization constant tends to hurt model performance, even not using a regularization constant ($\alpha = 0$) hurts precision compared to using a small regularization constant.

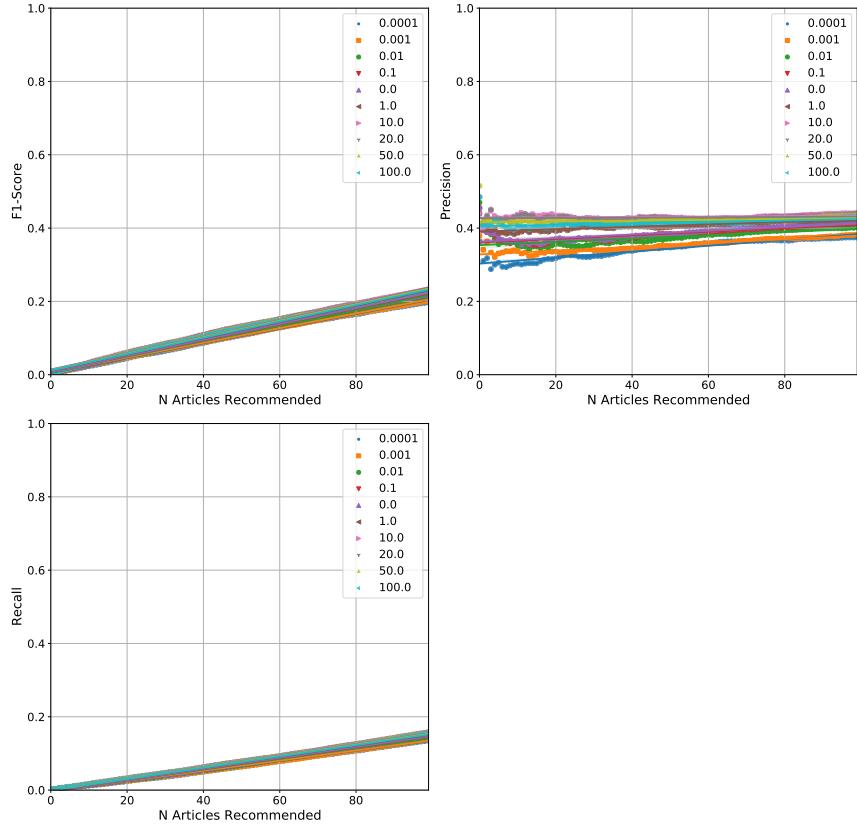
Regularization Constant vs Model Performance (Metrics @K) | Cumulative Performance across all Cluster Pairs --> Homogen



8.1.2 Heterogeneous Users

For Heterogeneous Users on the other hand, the greater the regularization constant the better the model performs (as it can be seen across all 3 metrics), so limiting spurious correlations definitely does help in this scenario as words that overlap across instances are demoted in importance.

Regularization Constant vs Model Performance (Metrics @K) | Cumulative Performance across all Cluster Pairs --> Heterogeneous

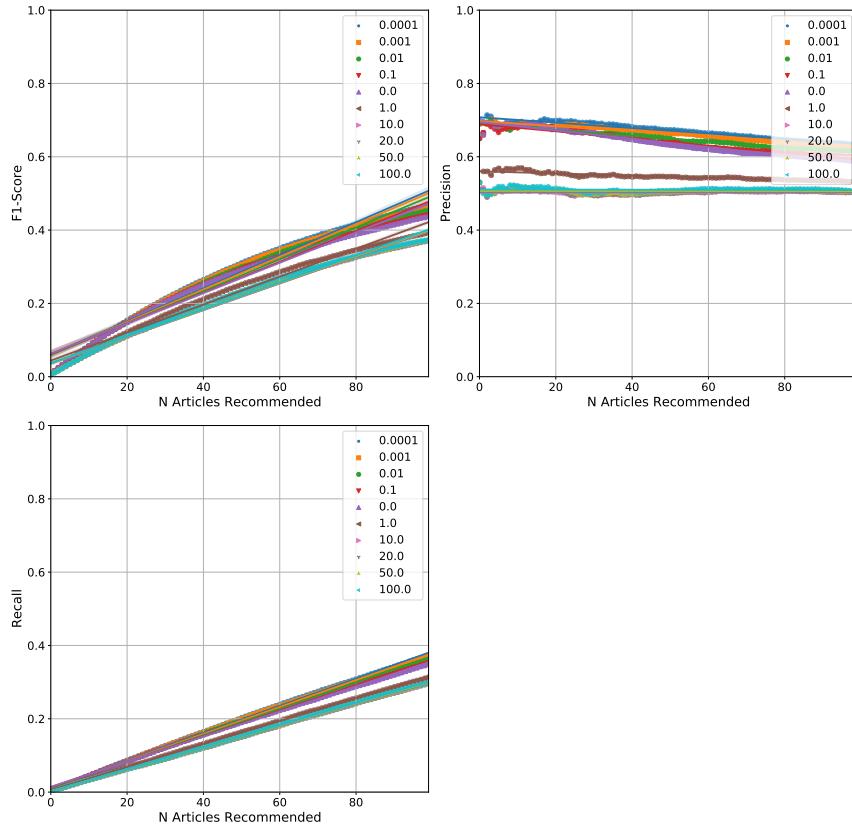


8.2 Glove

8.2.1 Homogeneous Users

For Homogeneous Users we see that a high regularization constant tends to keep the precision at a more constant rate compared to a smaller regularization constant but has lower values initially.

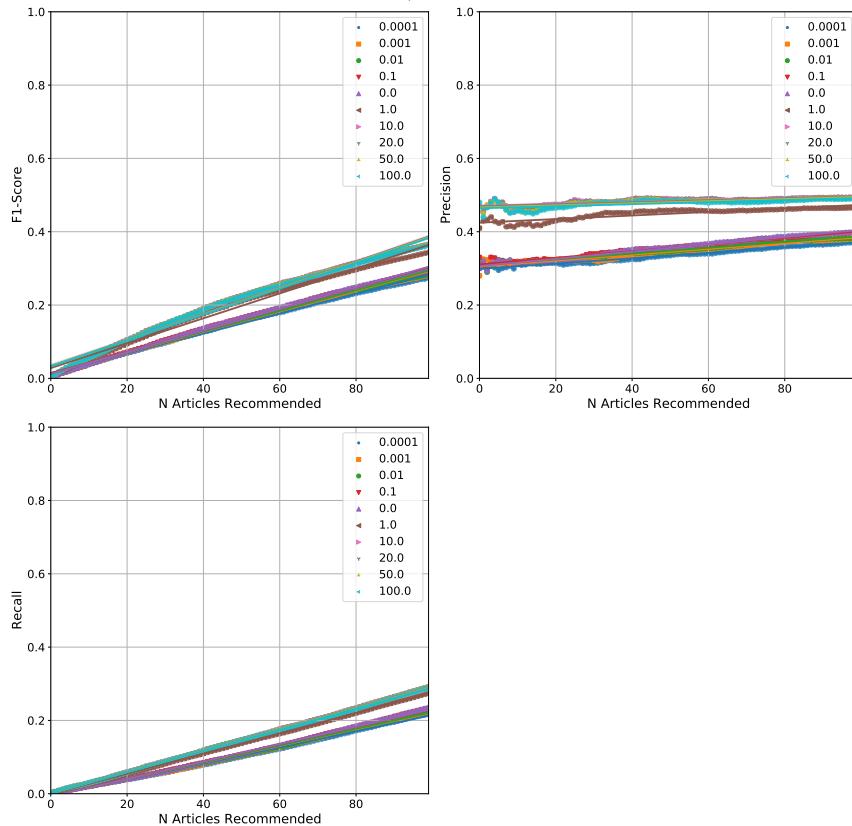
Regularization Constant vs Model Performance (Metrics @K) | Cumulative Performance across all Cluster Pairs --> Homogeneous



8.2.2 Heterogeneous Users

For Heterogeneous Users on the other hand, higher regularization constant values improves precision over lower regularization constant values and the precision seems to be increase at a more constant rate compared to lower values.

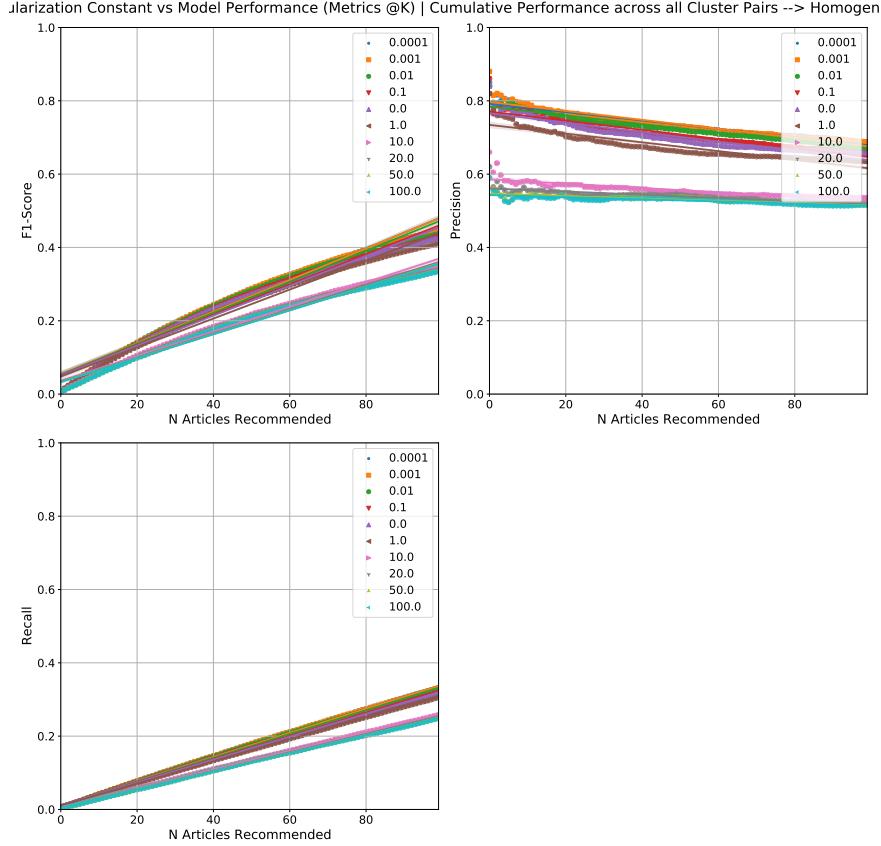
Regularization Constant vs Model Performance (Metrics @K) | Cumulative Performance across all Cluster Pairs --> Heterogeneous



8.3 BERT

8.3.1 Homogeneous Users

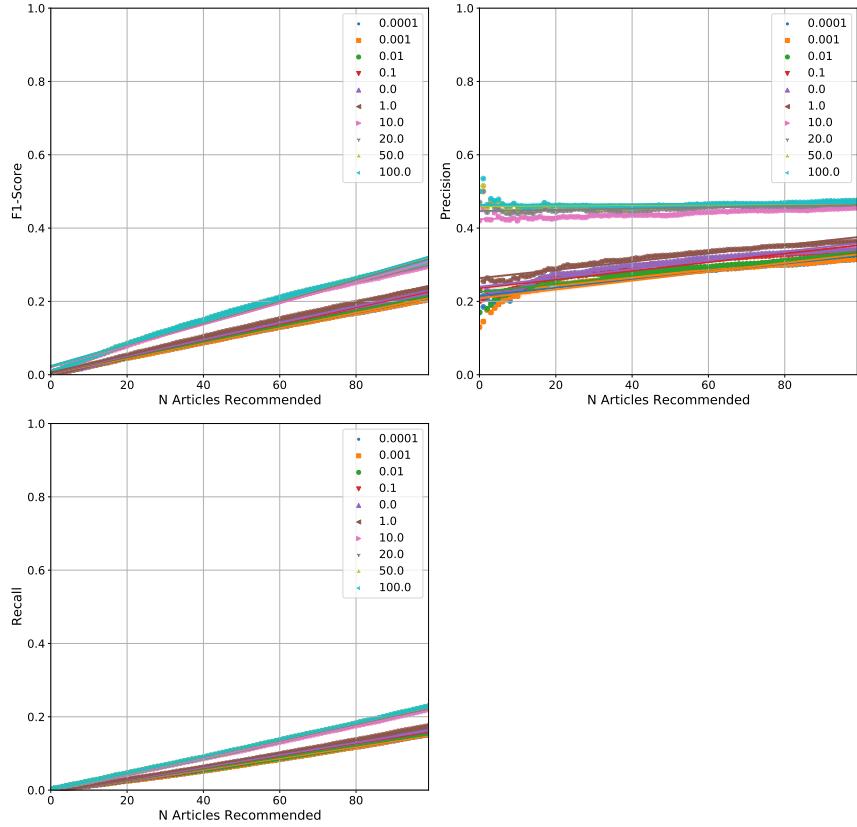
For Homogeneous Users we see that lower regularization values have high precision compared to higher constant values and we also observe that the precision for the initial interactions are much higher compared to when using TFIDF and GLOVE representations



8.3.2 Heterogeneous Users

For Heterogeneous Users on the other hand, higher regularization constant values improves precision over lower regularization constant values and the precision seems to be increase at a more constant rate compared to lower values. We also observe that both GLOVE and TFIDF Representations tend to have a higher precision compared to BERT for lower regularization values.

Regularization Constant vs Model Performance (Metrics @K) | Cumulative Performance across all Cluster Pairs --> Heterogeneous



8.4 Summary

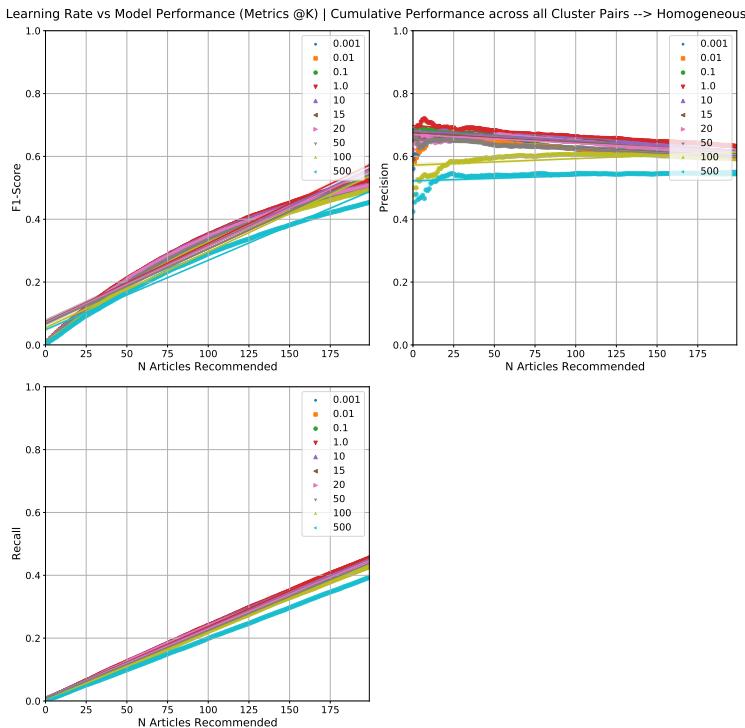
- TF-IDF is not affected greatly by using regularization constants for the Homogeneous users compared to GLOVE and BERT representations
- BERT representations have higher initial precision for Homogeneous users compared to TFIDF and GLOVE with low regularization constant values
- GLOVE representations with high regularization constant values seems to be performing the best for Heterogeneous Users, followed by BERT (with only high regularization values)

9 Baseline 5: Learning Rate vs Model Performance

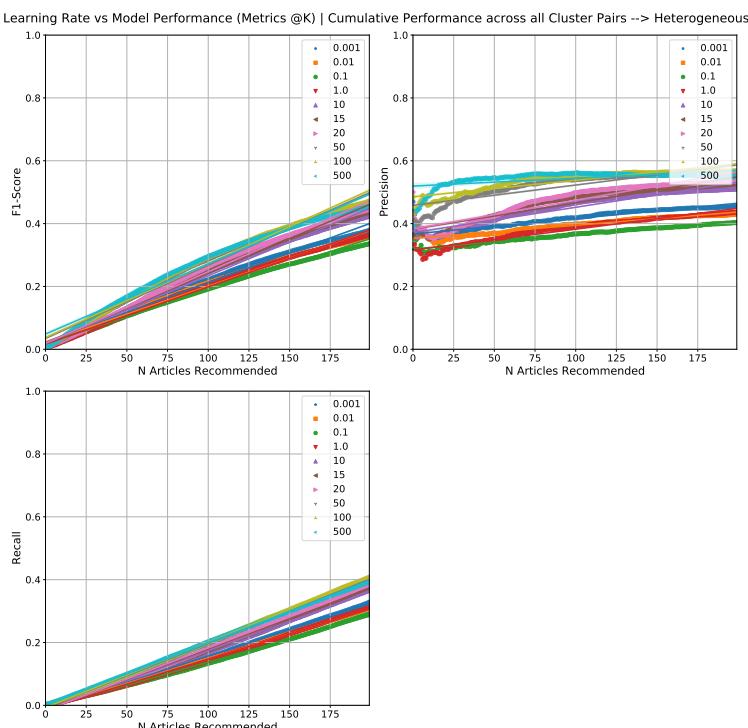
We want to measure the effect of using different learning rates to see which leads to better convergence when max-iterations are set to 1000.

9.1 TF-IDF

9.1.1 Homogeneous Users :

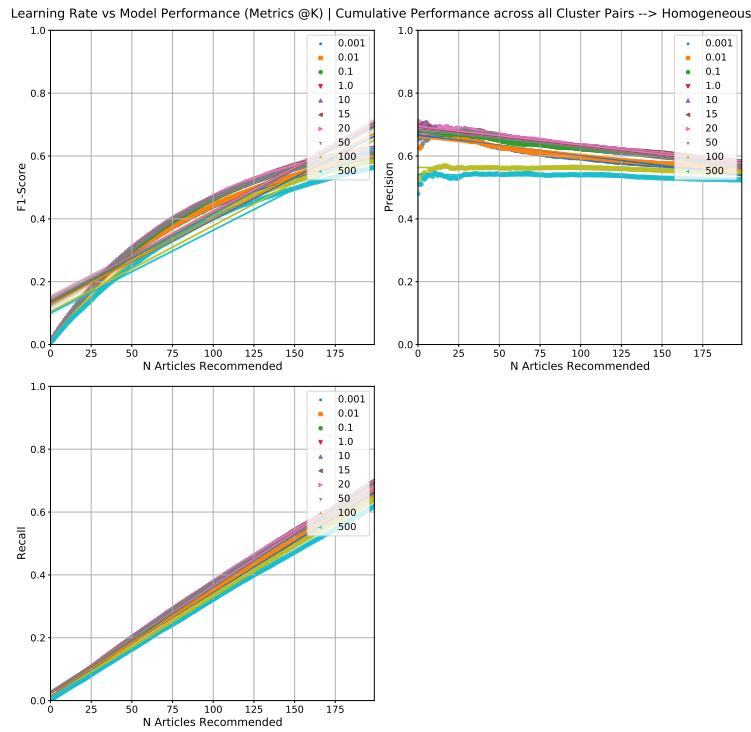


9.1.2 Heterogeneous Users :

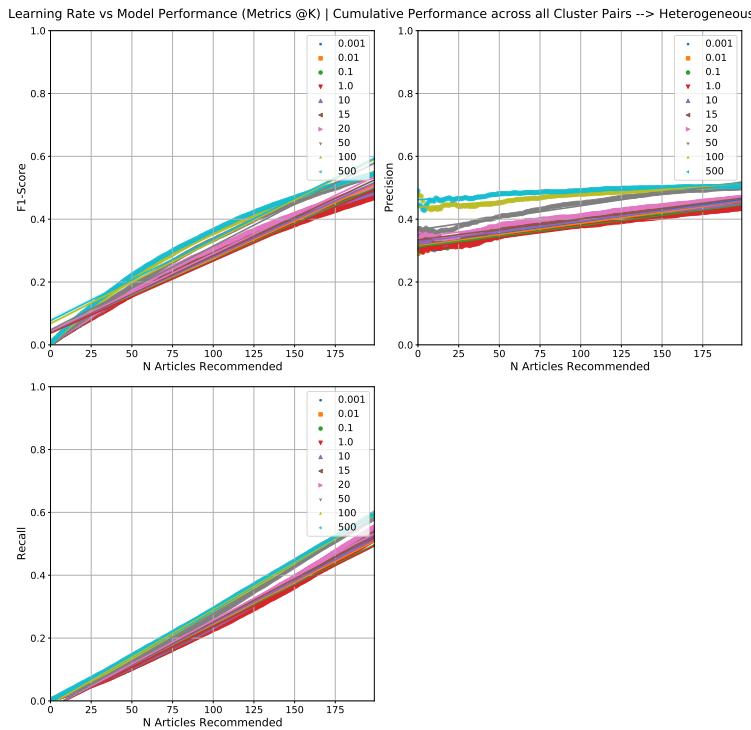


9.2 Glove

9.2.1 Homogeneous Users :

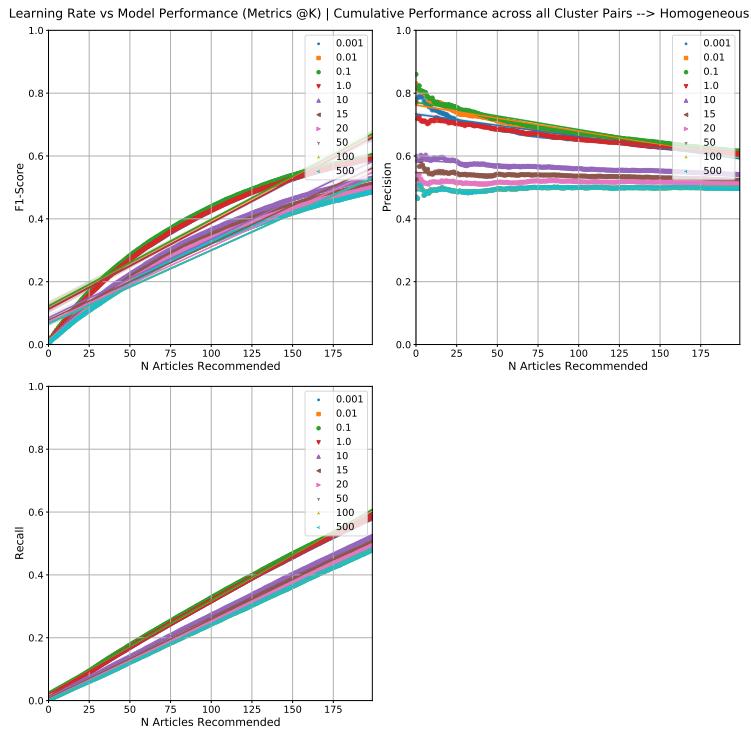


9.2.2 Heterogeneous Users :

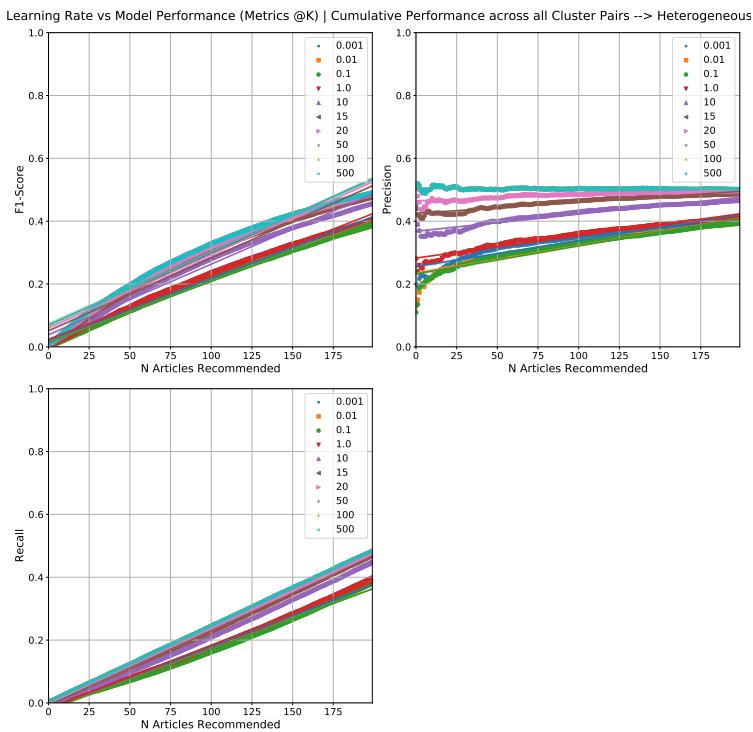


9.3 BERT

9.3.1 Homogeneous Users :



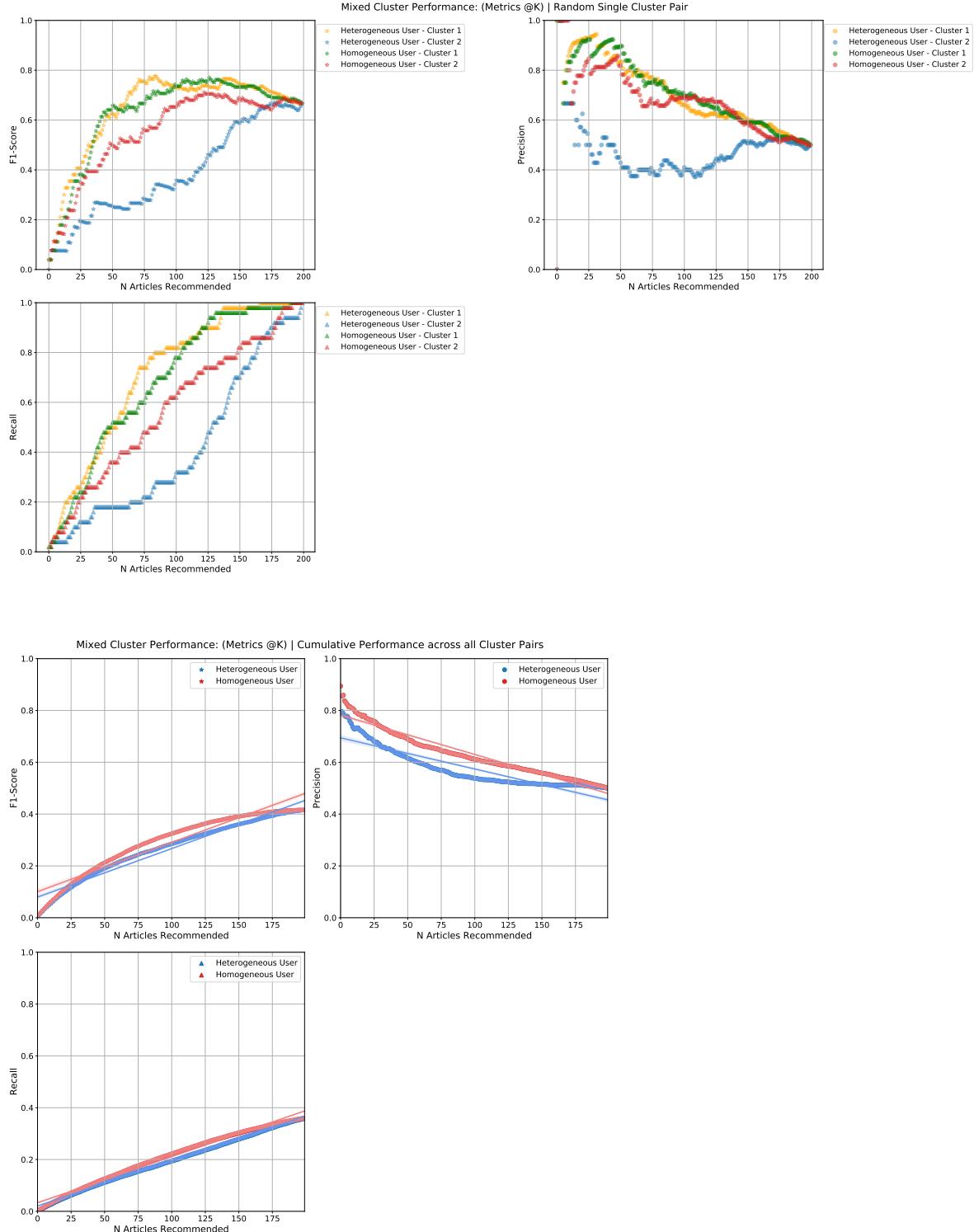
9.3.2 Heterogeneous Users :

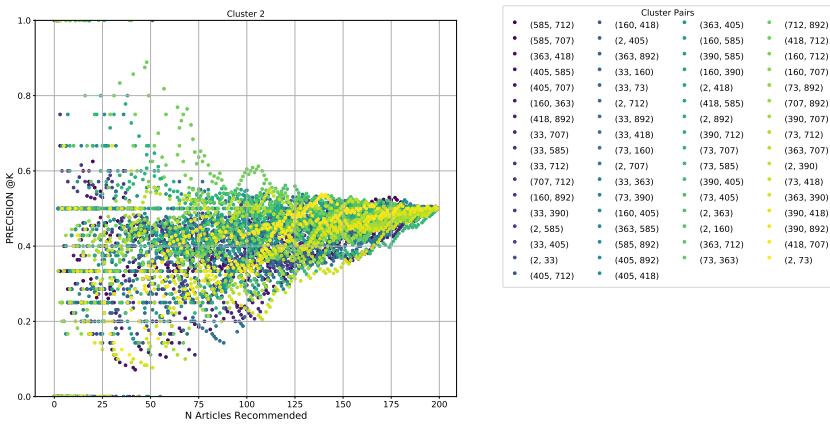
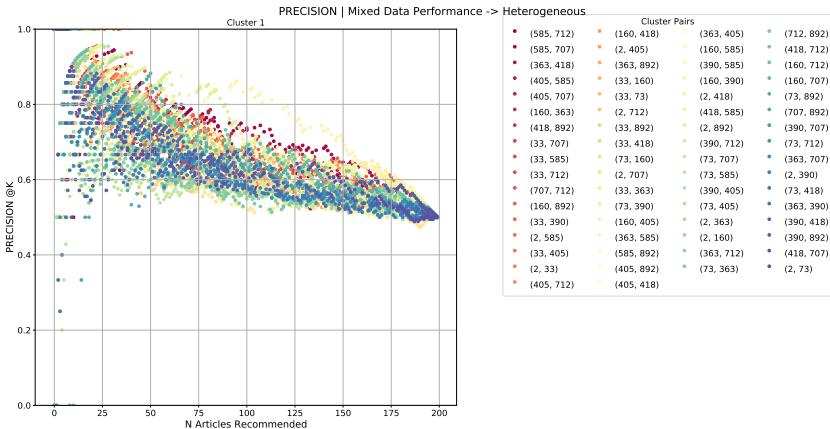
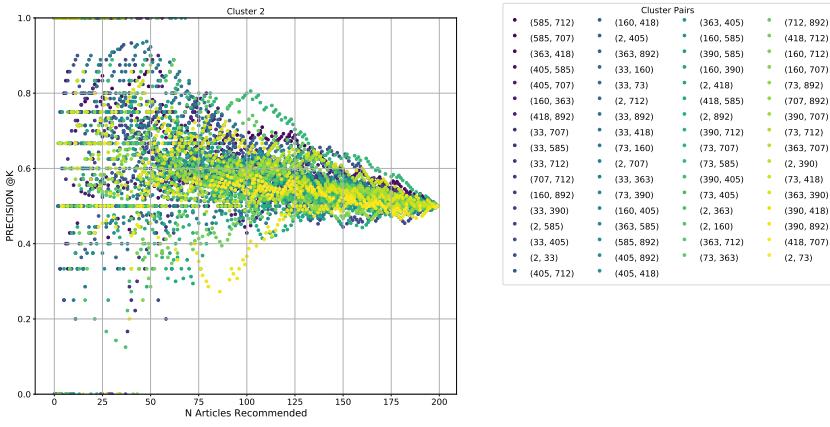
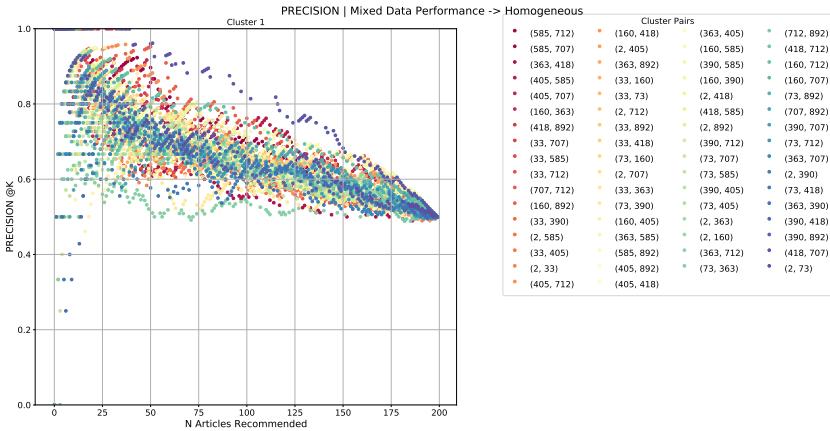


10 Baseline 6: Model Performance on mixed Set (Cluster 1 + Cluster 2)

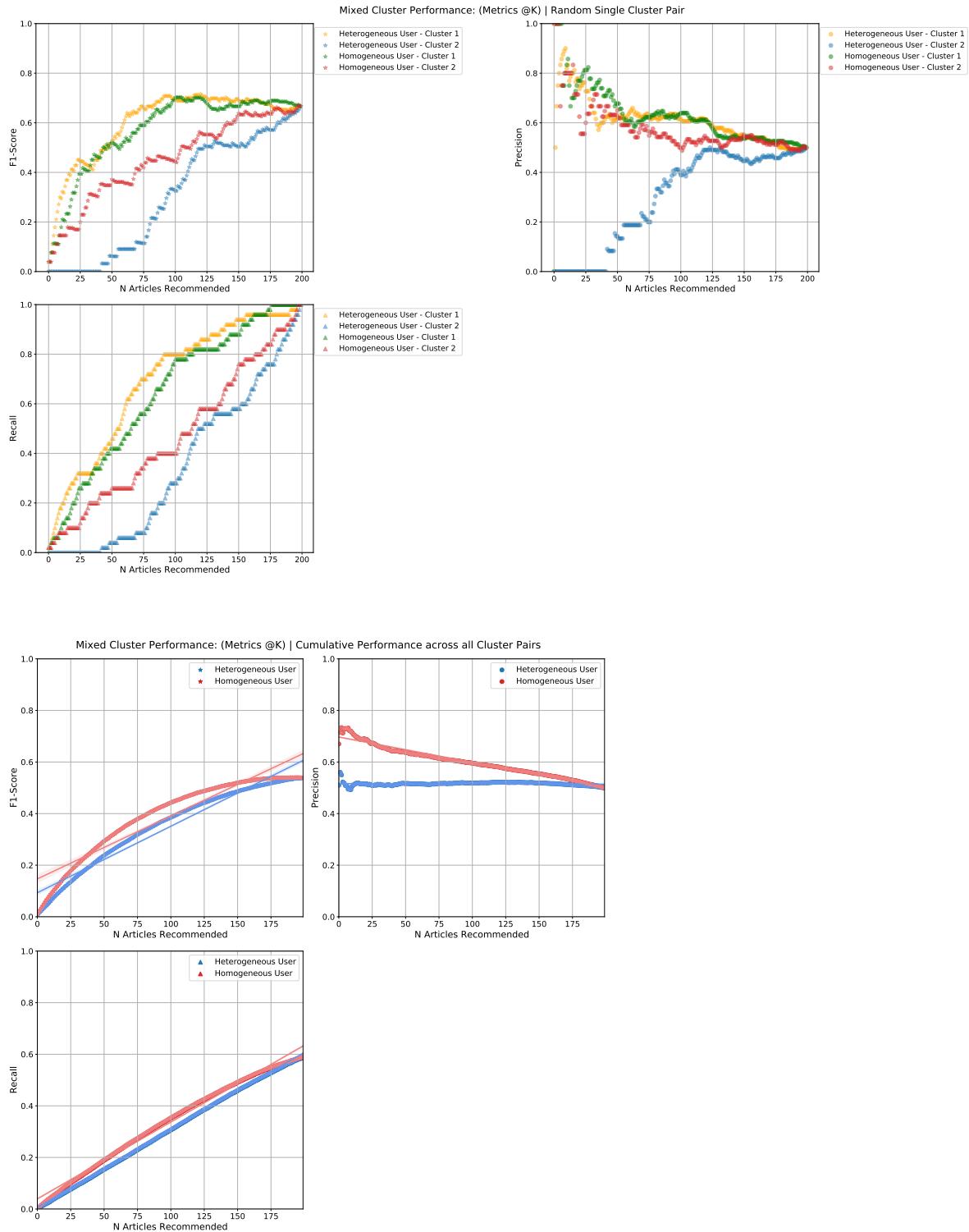
We want to measure how well the model would perform when an article from an unseen topic arrives along with articles from the original topic the classifier was trained on

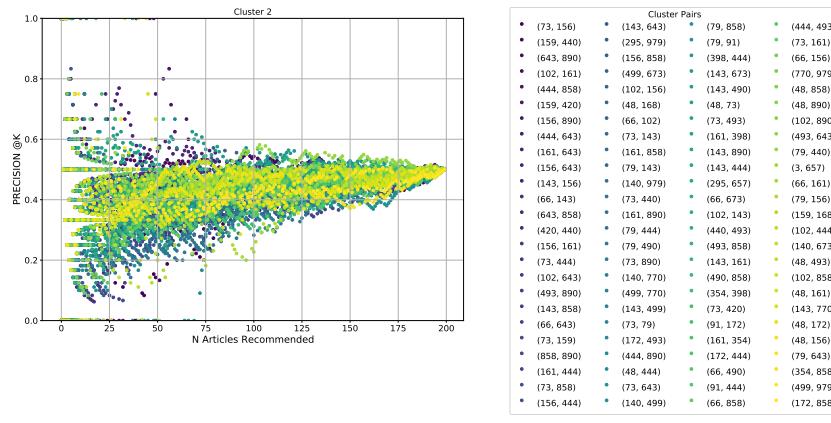
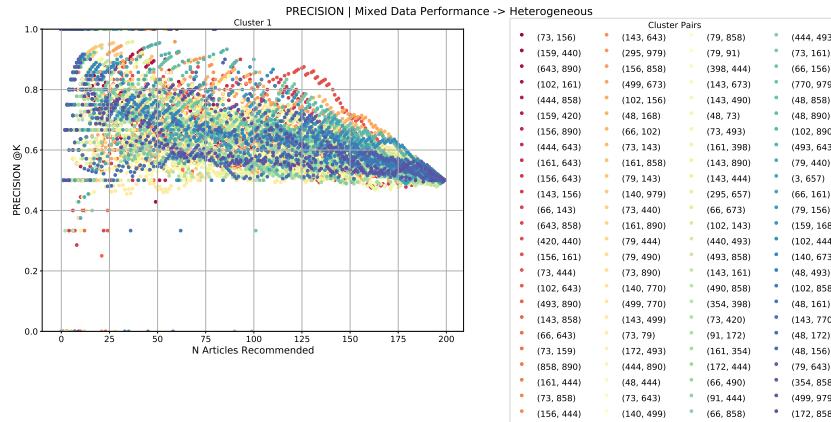
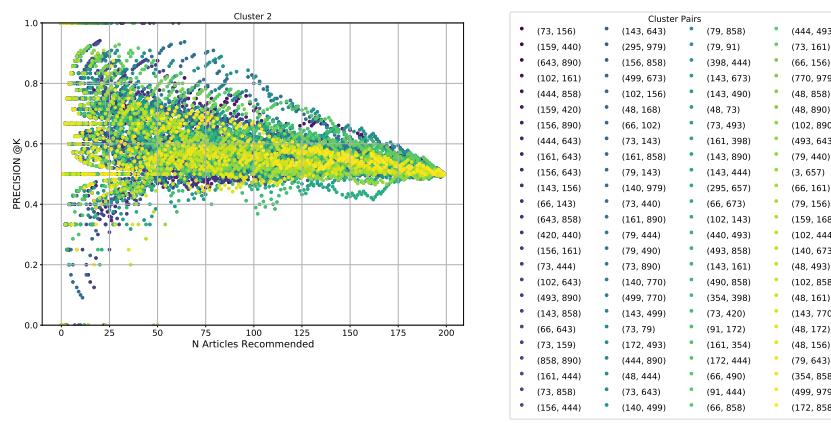
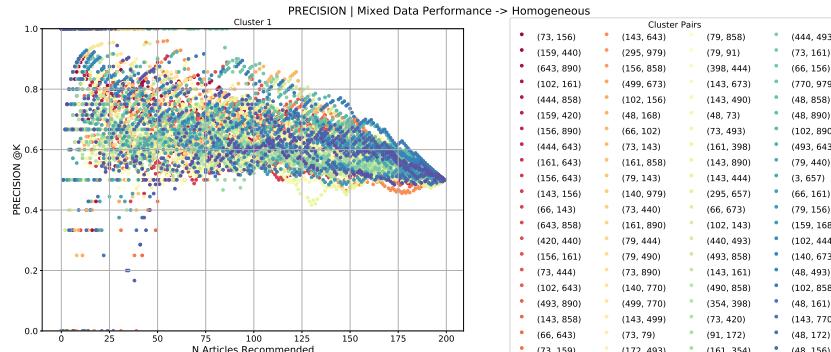
10.1 TF-IDF



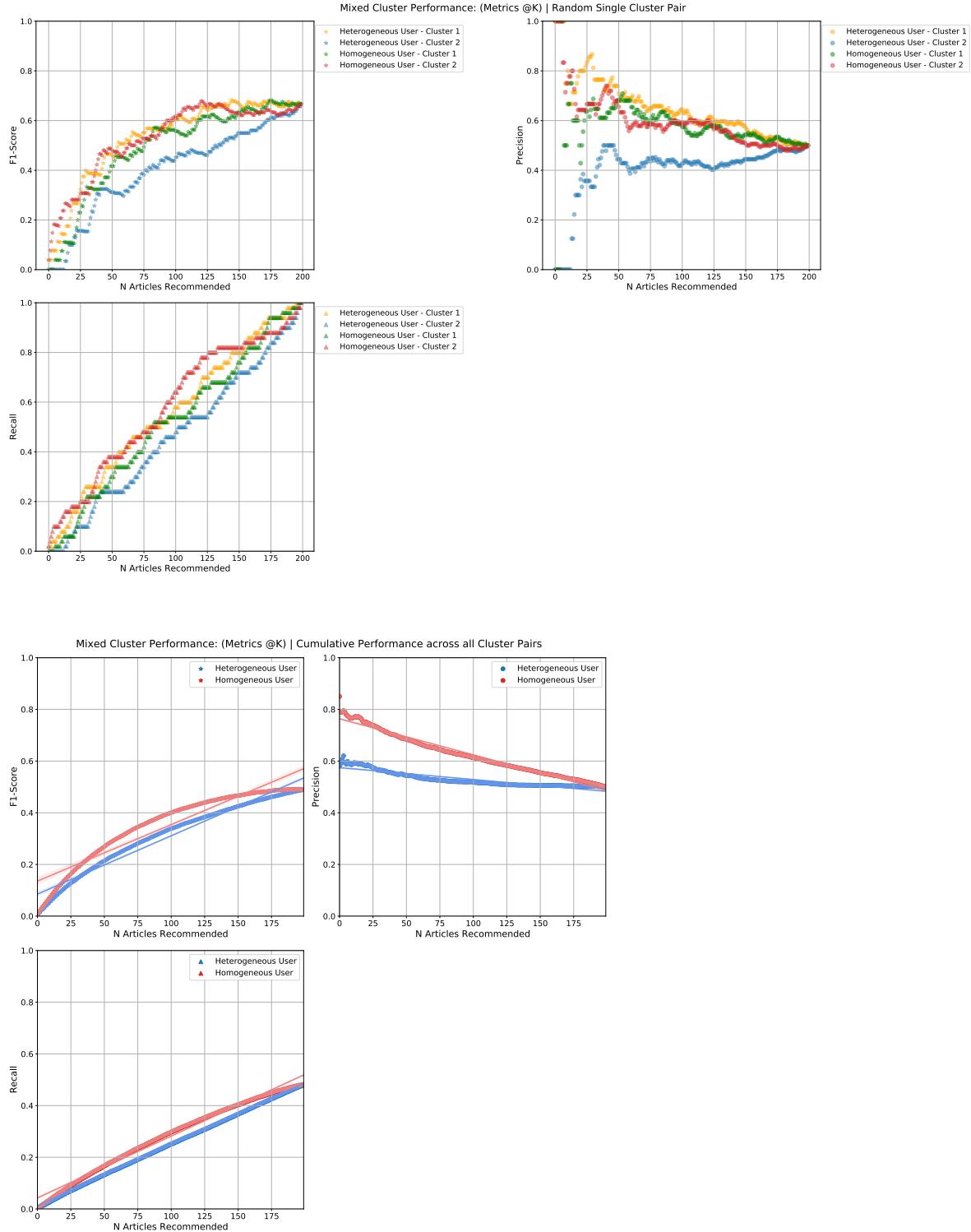


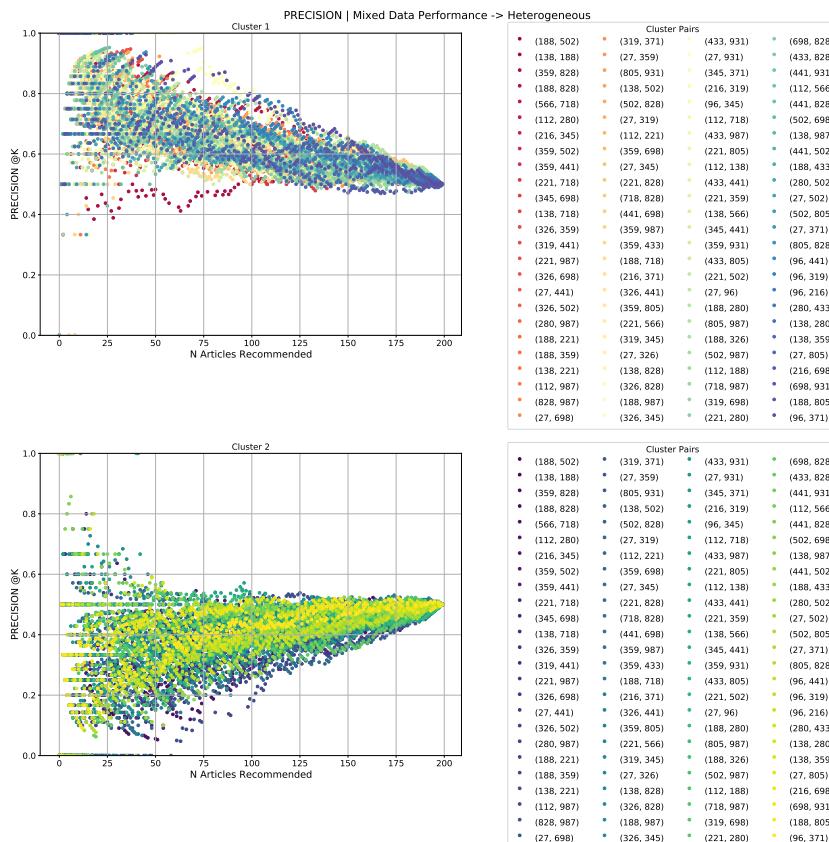
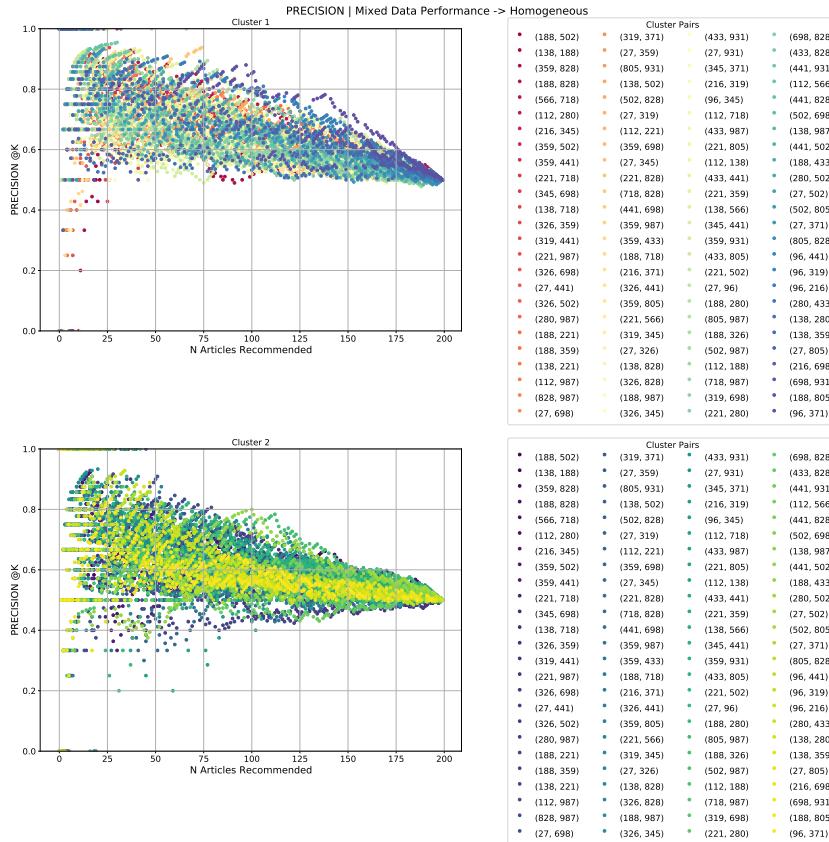
10.2 Glove





10.3 BERT

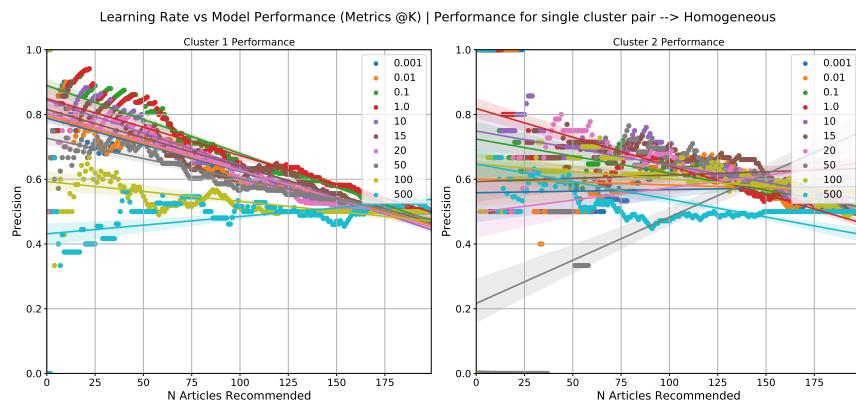


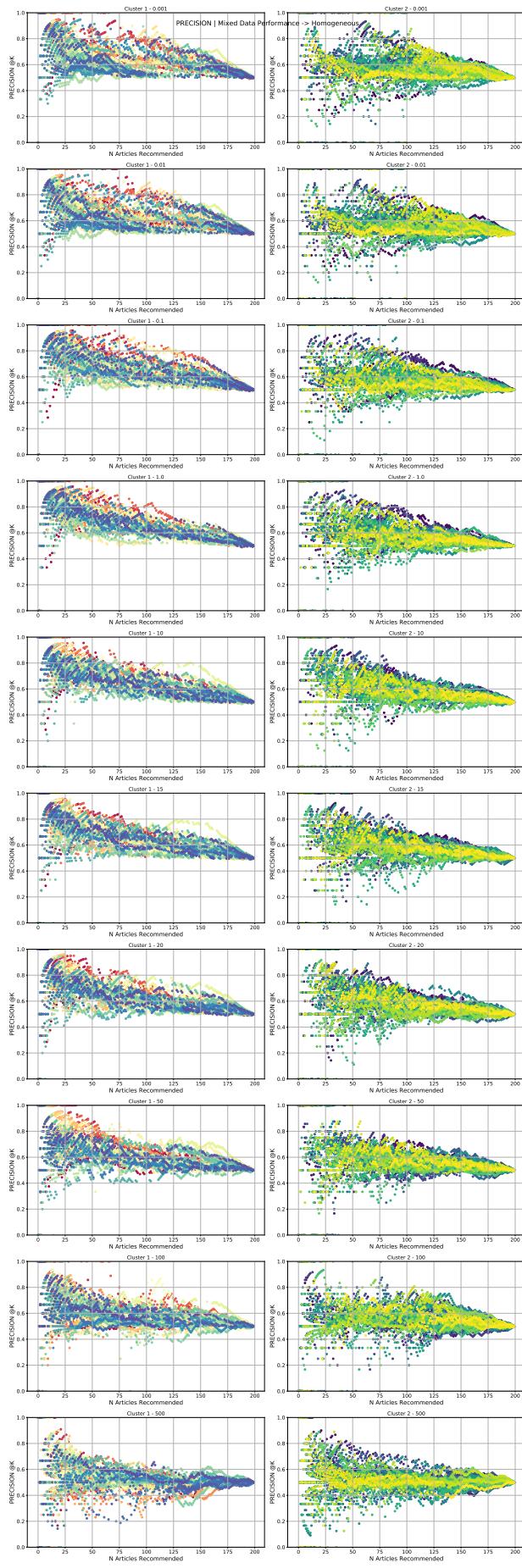


11 Baseline 7: Learning Rate Variation for Mixed Cluster

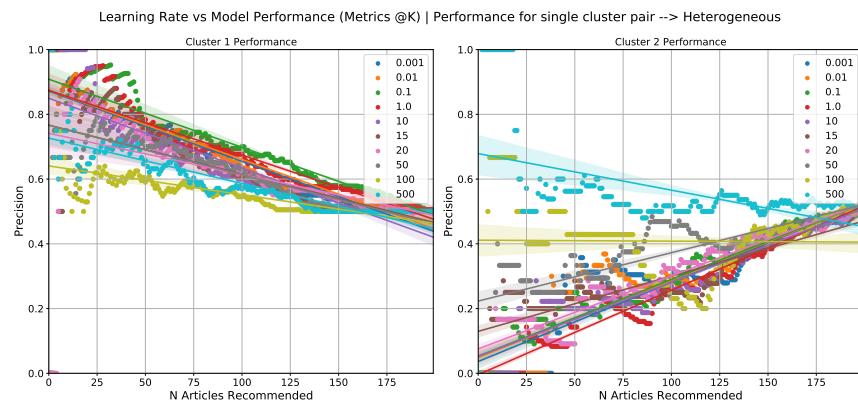
11.1 TF-IDF

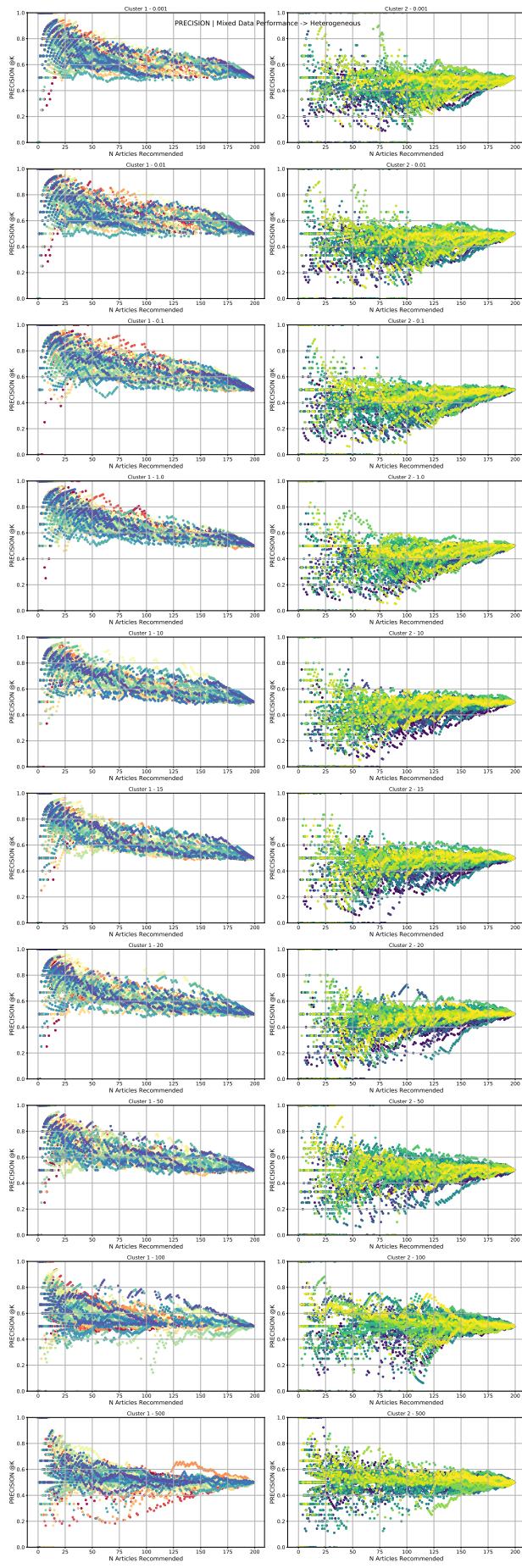
11.1.1 Homogeneous





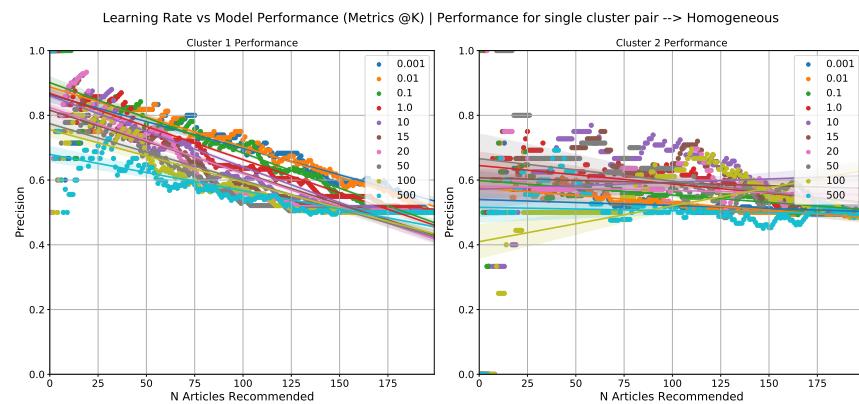
11.1.2 Heterogeneous

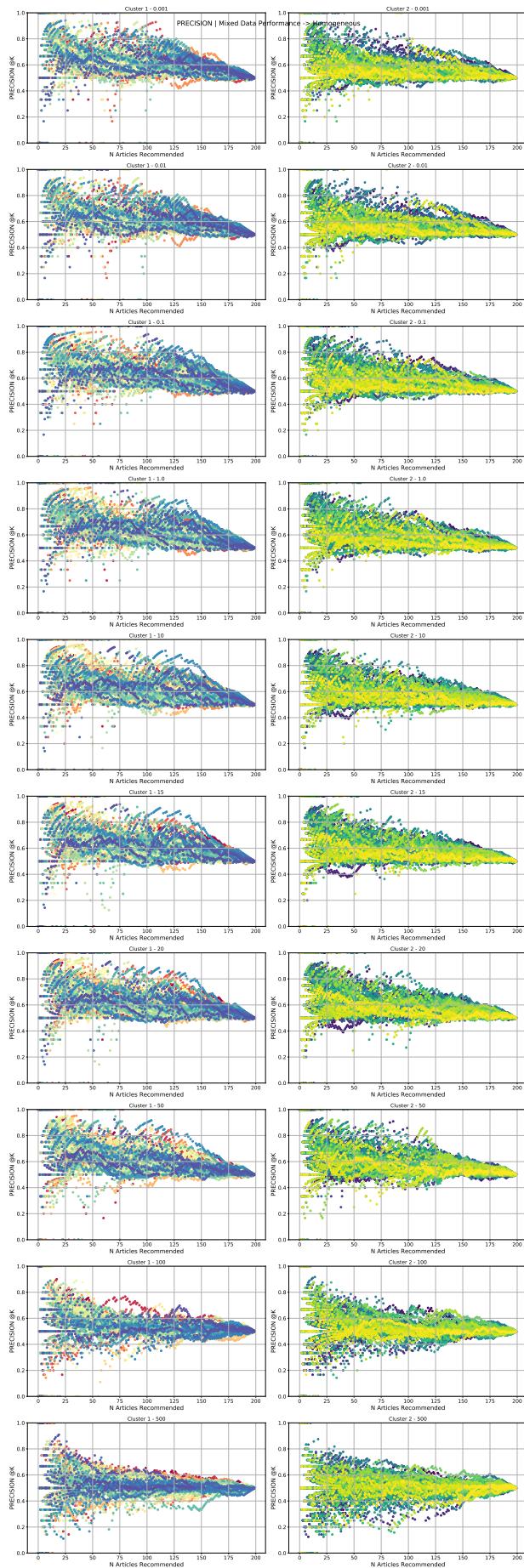




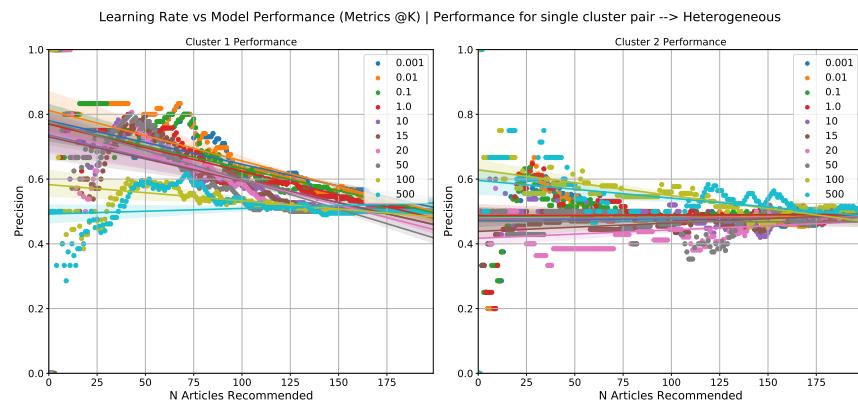
11.2 Glove

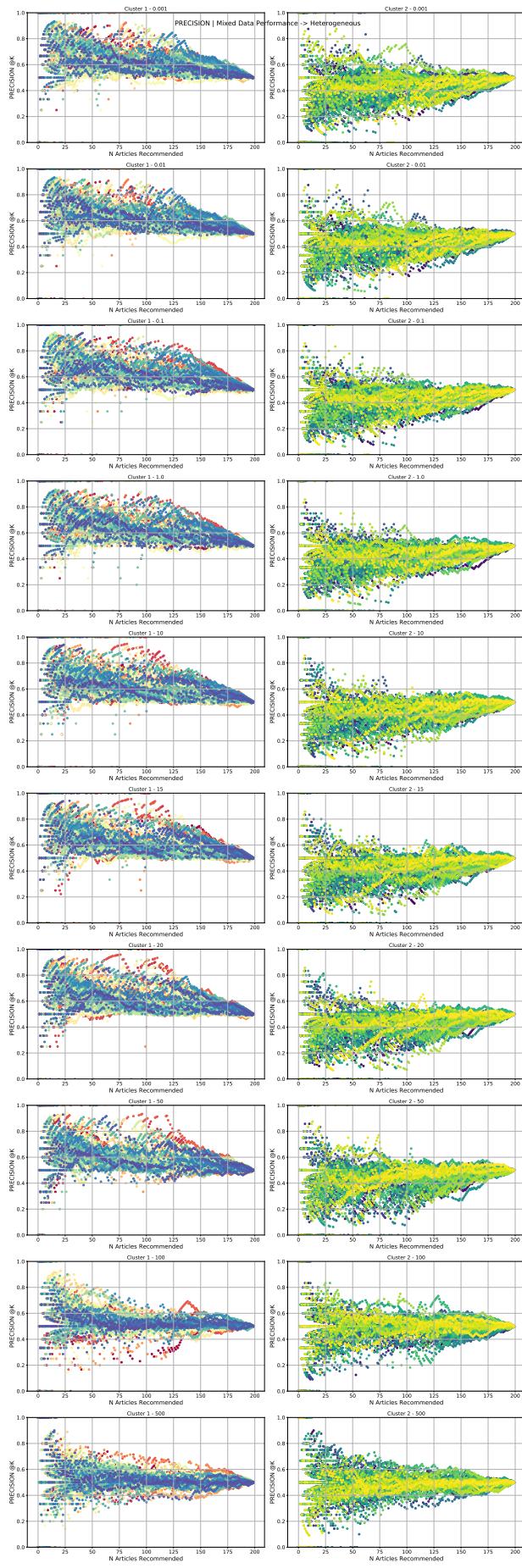
11.2.1 Homogeneous





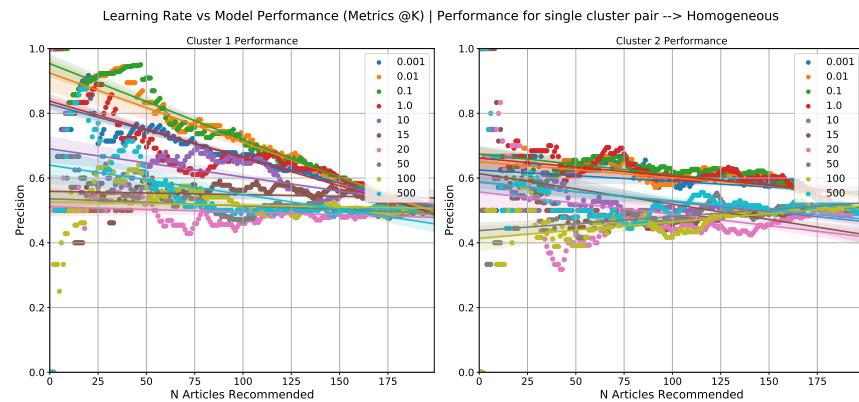
11.2.2 Heterogeneous

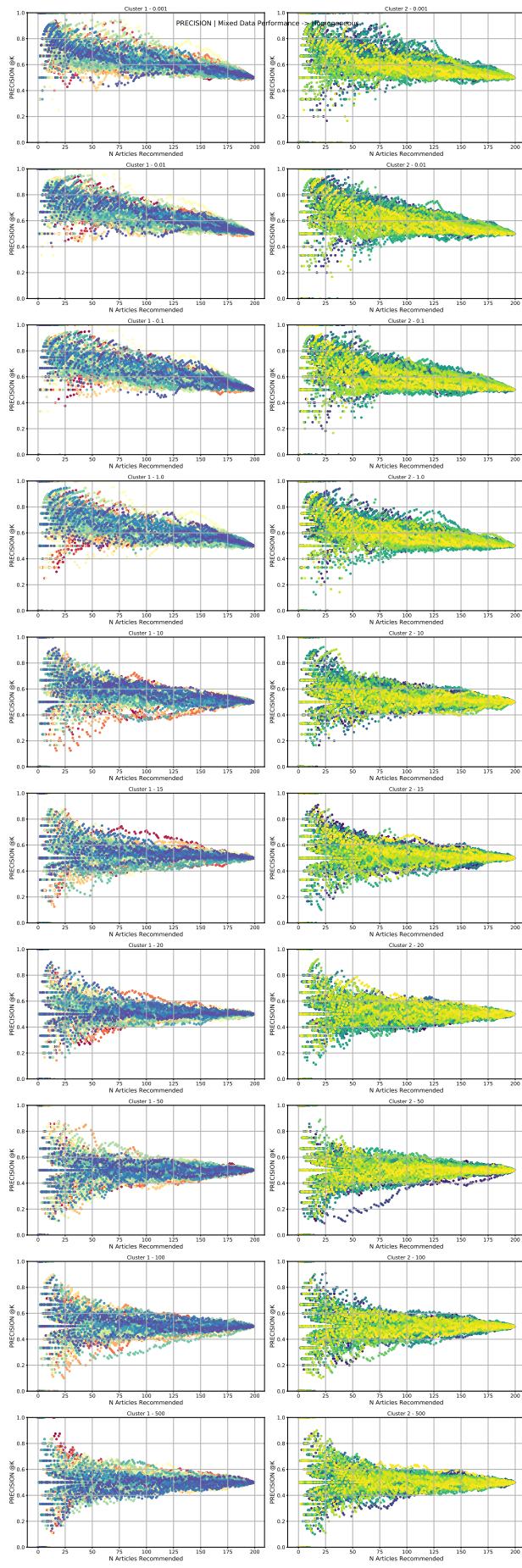




11.3 BERT

11.3.1 Homogeneous





11.3.2 Heterogeneous

