```python
# computational imports
import numpy as np
import pandas as pd
from ast import literal_eval
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import linear_kernel
import nltk
from nltk.tokenize import sent_tokenize
from nltk import word_tokenize
nltk.download('averaged_perceptron_tagger')
from sklearn.feature_extraction import text
from nltk.stem import WordNetLemmatizer
from nltk.corpus import wordnet as wn
import string
# plotting imports
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style("darkgrid")
from scipy.spatial import distance
# for reading
import urllib.request
# display imports
from IPython.display import display, IFrame
from IPython.core.display import HTML
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from sklearn.metrics.pairwise import cosine_similarity
from surprise import Reader, Dataset, KNNBasic, NormalPredictor,Baseline(
from surprise import SVD, SVDpp, NMF, SlopeOne, CoClustering
from surprise.model_selection import cross_validate
from surprise.model_selection import GridSearchCV
from surprise import accuracy
import random
from ast import literal_eval
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import linear_kernel
import itertools

# plotting
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style("darkgrid")
plt.style.use('ggplot')

# for reading files
import urllib.request
from nltk.stem import WordNetLemmatizer
from nltk.corpus import wordnet as wn
from nltk.tokenize import sent_tokenize

# display imports
```

```
57  from IPython.display import display, IFrame
58  from IPython.core.display import HTML
59  import surprise
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]     C:\Users\Dawit\AppData\Roaming\nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]       date!
```

# Load data

In [2]:
```
1  #read in the file
2  df = pd.read_csv('movies_metadata_clean.csv')
3  df = df.drop_duplicates(subset='title', keep="first", inplace=False)
4  df= df[0:3000]
5  print(f'The shape of the dataframe is {df.shape}')
6
7  # Convince Python that this column should be treated like a list, not a
8  df['genres'] = df['genres'].apply(literal_eval)
9  df[1:2]
```

The shape of the dataframe is (3000, 10)

Out[2]:

| | id | title | budget | genres | overview | revenue | runtime | vote_average | vote_ |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 8844 | Jumanji | 65000000.0 | [Adventure, Fantasy, Family] | When siblings Judy and Peter discover an encha... | 262797249.0 | 104.0 | 6.9 | |

In [3]:
```
1  movies = pd.DataFrame({'movie_id': df["id"],
2                         'title':df["title"],
3                         'genres': df["genres"]
4                         })
```

# Generate users and ratings for all movies in the dataset

In [4]: ▶|
```python
 1  n_users = 401
 2  n_movies = 3000
 3  #generate a rating for each user/movie combination
 4  ratings = pd.DataFrame(np.array(np.meshgrid([userid for userid in range(
 5  np.random.seed(12)
 6  randratings = np.random.randint(1,6, ratings.shape[0])
 7  ratings['rating'] = randratings
 8  ratings["user_id"] = ratings["user_id"].astype(str)
 9  ratings["movie_id"] = ratings["movie_id"].astype(str)
10  ratings.dtypes, ratings.shape
```

Out[4]: (user_id     object
         movie_id    object
         rating       int32
         dtype: object,
         (1200000, 3))

# Compute Similarity among movies in the dataset

*Lemma Tokenizer function to extract root words from text*

```
1  def get_wordnet_pos(word, pretagged = False):
2      """Map POS tag to first character lemmatize() accepts"""
3      if pretagged:
4          tag = word[1].upper()
5      else:
6          tag = nltk.pos_tag([word])[0][1][0].upper()
7      tag_dict = {"J": wn.ADJ,
8                  "N": wn.NOUN,
9                  "V": wn.VERB,
10                 "R": wn.ADV}
11
12     return tag_dict.get(tag, wn.NOUN)
13
14
15 class LemmaTokenizer(object):
16     def __init__(self):
17         self.wnl = WordNetLemmatizer()
18     def __call__(self, articles):
19         sents = sent_tokenize(articles)
20         sent_pos = [nltk.pos_tag(word_tokenize(s)) for s in sents]
21         pos = [item for sublist in sent_pos for item in sublist]
22         lems = [self.wnl.lemmatize(t[0], get_wordnet_pos(t, True)) for t
23         lems_clean = [''.join(c for c in s if c not in string.punctuatio
24         return lems_clean
25
26 lemmatizer = WordNetLemmatizer()
27 lemmatized_stop_words = [lemmatizer.lemmatize(w) for w in text.ENGLISH_S
28 lemmatized_stop_words.extend(['ve','nt','ca','wo','ll'])
29
```

## Obtain *TF - IDF Scores for all movies based on the *overview column*

### 1. Intialize vetorizer

### 2. Remove stop words in the vector

### 3. Costruct TD-IDF matrix for all movies

```
1  df = df.copy()
2  tfidf = TfidfVectorizer(tokenizer=LemmaTokenizer(), lowercase=True, stop_
3  vec3 = CountVectorizer(tokenizer=LemmaTokenizer(), lowercase=True, stop_v
4  df['overview'] = df['overview'].fillna('')
5  tfidf_matrix = tfidf.fit_transform(df['overview'])
6  tfidf_matrix.shape
```

Out[8]: (3000, 100)

```
1  feature_names = tfidf.get_feature_names()
2  corpus_index = df['title']
3  pd.DataFrame(tfidf_matrix.todense(), index=corpus_index, columns=feature_
```

Out[9]:

| title | american | attempt | beautiful | begin | best | boy | brother | child | city | ... | v |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Jumanji | 0.510952 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... |
| Grumpier Old Men | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... |
| Waiting to Exhale | 0.610904 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... |
| Father of the Bride Part II | 0.504329 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... |

4 rows × 100 columns

### Compute Cosine Similarity using TD-IDF scores and store in a matrix
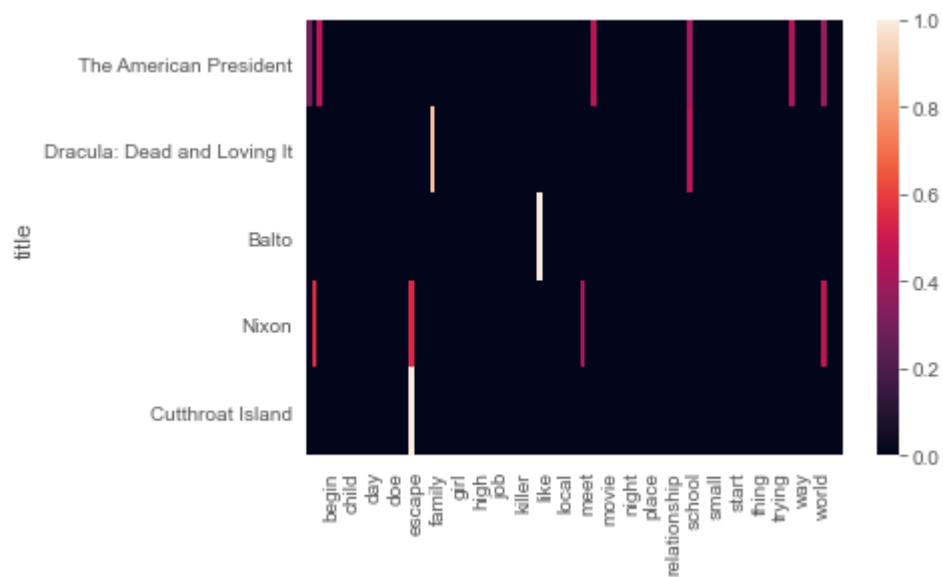
In [10]:

```
1  # Compute the cosine similarity matrix
2  sim_matrix = linear_kernel(tfidf_matrix, tfidf_matrix)
3  dpdf = pd.DataFrame(sim_matrix, columns=df['title'], index=df['title'])
4  dpdf.shape
```

Out[10]: (3000, 3000)

-------------------------------------------------------------------------

# Content Based Recommendation

### TD-IDF visulization

```
1 ax = sns.heatmap(pd.DataFrame(tfidf_matrix.todense()[10:15], index=df['t:
```



## Cosine similairty table

In [12]:

```
1 sample = dpdf.iloc[1:5,1:5]
2 cm = sns.color_palette("Blues", as_cmap=True)
3 sample.style.set_caption('Dot Product with most similar movies highlighte
4     .background_gradient(cmap=cm)
5 # sample
```

Out[12]:

Dot Product with most similar movies highlighted.

| title | Jumanji | Grumpier Old Men | Waiting to Exhale | Father of the Bride Part II |
|---|---|---|---|---|
| **title** | | | | |
| **Jumanji** | 1.000000 | 0.064428 | 0.312143 | 0.440489 |
| **Grumpier Old Men** | 0.064428 | 1.000000 | 0.000000 | 0.084790 |
| **Waiting to Exhale** | 0.312143 | 0.000000 | 1.000000 | 0.308097 |
| **Father of the Bride Part II** | 0.440489 | 0.084790 | 0.308097 | 1.000000 |

### Content- based recommendation inital function

```
In [123]:   1  def content_recommender(df, seed, seedC, simmatrix, top):
            2      indices = pd.Series(df.index, index=df[seedC]).drop_duplicates()
            3      idx = indices[seed]
            4      sim_s = list(enumerate(simmatrix[idx]))
            5      del sim_s[idx]
            6      sim_s = sorted(sim_s, key=lambda x: x[1], reverse=True)
            7      sim_s = sim_s[:top]
            8      movie_idx = [i[0] for i in sim_s]
            9      return df.iloc[movie_idx]
```

### Test function

```
In [124]:  ▶|   1  content_recommender(df, 'The Locusts', 'title', sim_matrix, 10)
```

Out[124]:

| | id | title | budget | genres | overview | revenue | runtime | vote_av |
|---|---|---|---|---|---|---|---|---|
| **2056** | 161795 | Déjà Vu | 0.0 | [Romance, Drama] | L.A. shop owner Dana and Englishman Sean meet ... | 0.0 | 117.0 | |
| **57** | 11010 | The Postman | 0.0 | [Comedy, Drama, Romance] | Simple Italian postman learns to love poetry w... | 0.0 | 108.0 | |
| **1372** | 2892 | Angel Baby | 0.0 | [Drama] | Two schizophrenics meet during therapy and fal... | 0.0 | 105.0 | |
| **447** | 25440 | Widows' Peak | 0.0 | [Comedy, Thriller, Mystery, Romance, Foreign] | Scandal and mystery reign following the arriva... | 0.0 | 101.0 | |
| **90** | 9095 | Mary Reilly | 47000000.0 | [Drama, Horror, Thriller, Romance] | A housemaid falls in love with Dr. Jekyll and ... | 12379402.0 | 104.0 | |
| **483** | 1413 | M. Butterfly | 0.0 | [Drama, Romance] | In 1960s China, French diplomat Rene Gallimard... | 1499795.0 | 101.0 | |
| **566** | 95743 | Foreign Student | 0.0 | [Drama, Romance] | A French football playing exchange student fal... | 0.0 | 90.0 | |
| **759** | 32872 | Til There Was You | 1000000.0 | [Comedy, Romance] | Two strangers, whose paths are always crossing... | 0.0 | 113.0 | |
| **2923** | 2039 | Moonstruck | 0.0 | [Comedy, Drama, Romance] | Cher is devastatingly funny, sinuous and beaut... | 80640528.0 | 102.0 | |
| **1812** | 65203 | The Broadway Melody | 379000.0 | [Drama, Music, Romance] | Harriet and Queenie Mahoney, a vaudeville act,... | 4358000.0 | 100.0 | |

*Function for conent based recommendation for N users*

```
1  # list_of_movies = df["title"]
2  def content_recommender2(n_users,n_items,df):
3      list_of_movies = df["title"][1:n_users+1]
4      list_of_recommendations=[]
5      for movie in list_of_movies:
6          list_of_recommendations.append(list(content_recommender(df, str(
7      return list_of_recommendations
8
```

*Make recommendations for N users*

In [16]:

```
1  content_list_of_recommendations = content_recommender2(400,10,df)
```

----------------------------------------------------------------------------------------

# Collaborative - Filtering recommendation system using KNN

*KNN- Model setup*

In [17]: 

```python
 1  our_seed = 14
 2
 3  #Define a Reader
 4  reader = Reader(rating_scale=(1,5)) # defaults to (0,5)
 5
 6  #Create the dataset
 7  data = Dataset.load_from_df(ratings, reader)
 8
 9  #Define the algorithm object
10  knn = KNNBasic(k= 4, n_jobs=-1, verbose=False) #the default for k is 40,
11
12  random.seed(our_seed)
13  np.random.seed(our_seed)
14
15  #cross validation
16  knn_cv = cross_validate(knn, data, measures=['RMSE'], cv=5, verbose=True
17  print(knn_cv)
18
19  #extract RMSE
20  knn_RMSE = np.mean(knn_cv['test_rmse'])
21  print(f'\nThe RMSE across five folds was {knn_RMSE}')
```

Evaluating RMSE of algorithm KNNBasic on 5 split(s).

|                | Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 | Mean   | Std    |
|----------------|--------|--------|--------|--------|--------|--------|--------|
| RMSE (testset) | 1.5792 | 1.5810 | 1.5807 | 1.5799 | 1.5785 | 1.5799 | 0.0009 |
| Fit time       | 38.81  | 41.77  | 21.38  | 19.10  | 22.36  | 28.68  | 9.58   |
| Test time      | 125.70 | 94.71  | 69.75  | 71.24  | 73.90  | 87.06  | 21.33  |

```
{'test_rmse': array([1.57924483, 1.58103116, 1.58071122, 1.5799147 , 1.5785
2976]), 'fit_time': (38.809500217437744, 41.77214026451111, 21.377328634262
085, 19.097674131393433, 22.36443567276001), 'test_time': (125.703440904617
31, 94.70789527893066, 69.74647951126099, 71.23967361450195, 73.89880084991
455)}
```

The RMSE across five folds was 1.5798863348788867

### Grid search for hyperparameter tuning of the KNN model

In [18]: 

```python
 1  data.build_full_trainset()
```

Out[18]:  <surprise.trainset.Trainset at 0x197a4f89c48>

```python
In [19]:  1  #Define a Reader object
          2  #The Reader object helps in parsing the file or dataframe containing rati
          3  reader = Reader(rating_scale=(1,5)) # defaults to (0,5)
          4
          5
          6  #Create the dataset
          7  data = Dataset.load_from_df(ratings, reader)
          8  raw_ratings = data.raw_ratings
          9
         10  # shuffle ratings
         11  random.seed(our_seed)
         12  np.random.seed(our_seed)
         13  random.shuffle(raw_ratings)
         14
         15  #A = 90% of the data, B = 10% of the data
         16  threshold = int(.9 * len(raw_ratings))
         17  A_raw_ratings = raw_ratings[:threshold]
         18  B_raw_ratings = raw_ratings[threshold:]
         19
         20  data.raw_ratings = A_raw_ratings   # data is now the set A
         21
         22  # Select your best algo
         23  print('Grid Search...')
         24  param_grid = {'k': [3,5], 'min_k': [1,3]} #this will all combinations of
         25  grid_search = GridSearchCV(KNNBasic, param_grid, measures=['rmse'], cv=3
         26  grid_search.fit(data)
         27  knn_gs_algo = grid_search.best_estimator['rmse']
         28
         29  # retrain on the whole set A
         30  trainset = data.build_full_trainset()
         31  knn_gs_algo.fit(trainset)
         32
         33  # Compute biased accuracy on A
         34  predictions = knn_gs_algo.test(trainset.build_testset())
         35  print(f'Biased accuracy on A = {accuracy.rmse(predictions)}')
         36
         37  # Compute unbiased accuracy on B
         38  testset = data.construct_testset(B_raw_ratings)   # testset is now the se
         39  predictions = knn_gs_algo.test(testset)
         40  print(f'Unbiased accuracy on B = {accuracy.rmse(predictions)}')
         41
```

```
Grid Search...
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
```

```
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
RMSE: 0.6917
Biased accuracy on A = 0.6916816445638049
RMSE: 1.5500
Unbiased accuracy on B = 1.549958370462392
```

In [20]: ▶|
```python
1  #we can see what our best parameters were
2  grid_search.best_params['rmse']
```

Out[20]: {'k': 5, 'min_k': 1}

## Retrain model using new parameters

In [21]: ▶|
```python
1   #set our seeds again
2   random.seed(our_seed)
3   np.random.seed(our_seed)
4
5   #reset the data.raw_ratings to 100%
6   data.raw_ratings = raw_ratings
7
8   #trainset
9   trainset = data.build_full_trainset()
10
11  #build the algorithm using the best parameters
12  knn_gs_algo = grid_search.best_estimator['rmse']
13
14  #fit to the data
15  knn_gs_algo.fit(trainset)
16
17  #predict user 1, movie 11
18  knn_gs_algo.predict("1","9626")
```

```
Computing the msd similarity matrix...
Done computing similarity matrix.
```

Out[21]: Prediction(uid='1', iid='9626', r_ui=None, est=1.787987652965251, details=
{'actual_k': 5, 'was_impossible': False})

## Test recommendation for user 1

```
In [22]:  ▶|    1  sample =  movies.copy()
                2  sample['est_rating'] = sample.apply(lambda x: knn_gs_algo.predict("1", x
                3  sample.sort_values('est_rating', ascending=False)[1:10]
```

Out[22]:

| | movie_id | title | genres | est_rating |
|---|---|---|---|---|
| **381** | 29444 | S.F.W. | [Comedy, Drama] | 4.887442 |
| **2387** | 18892 | Jawbreaker | [Comedy] | 4.887442 |
| **2525** | 15660 | Mommie Dearest | [Drama] | 4.887381 |
| **1390** | 14908 | McHale's Navy | [Action, Comedy, Romance] | 4.887381 |
| **385** | 315 | Faster, Pussycat! Kill! Kill! | [Action, Crime] | 4.887381 |
| **2958** | 11481 | Repulsion | [Drama, Horror, Thriller] | 4.887366 |
| **2667** | 28501 | The Pit and the Pendulum | [Fantasy, Horror, Drama] | 4.887366 |
| **1847** | 16619 | Ordinary People | [Drama] | 4.775084 |
| **197** | 79593 | The Tie That Binds | [Thriller] | 4.775084 |

### *Collaborative recommendation function*

```
In [23]:  ▶|    1  def collaborative_recommender(n_users,df):
                2      sample = df.copy()
                3      list_of_recommendations=[]
                4      for user in range(n_users+1):
                5          if user == 0:
                6              pass
                7          else:
                8              sample["estimated_rating"]= movies.apply(lambda x: knn_gs_alg
                9              list_of_recommendations.append(list(sample.sort_values('esti
               10      return list_of_recommendations
```

### *Make collaborative recommendations for N users*

```
In [24]:  ▶|    1  collaborative_list_of_recommendations = collaborative_recommender(400,mo
```

----------------------------------------------------------------------------------
--------

# Analyze recommendations from content-based and collaborative models above

*Function to compute similairty, diversity and total coverage for recommendations made using the two models*

In [25]:
```python
def uniqueCombinations(list_elements):
    l = list(itertools.combinations(list_elements, 2))
    s = set(l)
    return list(s)
```

In [26]:
```python
def compute_mean_sim_score(n_users,recommendations,sim_matrix,df):
    sim_scores_movies =[]
    score = 0
    N_unique_combinations =0
    for m in recommendations:
        for pair in uniqueCombinations(m):
            score = score + dpdf[pair[0]][pair[1]]
#           print(dpdf[pair[0]][pair[1]])
#           print(pair)
            N_unique_combinations += 1
        N_unique_combinations += 0
        sim_scores_movies.append(score)
        score = 0
    unique_combination_per_set =  N_unique_combinations /n_users
    mean_sim_scores_movies = [n/unique_combination_per_set for n in sim_
    total_recommendations = list(itertools.chain.from_iterable(recommenda
    number_Total_recommendations = len(set(total_recommendations))
    Total_movies = len(df['movie_id'])
    coverage = number_Total_recommendations /Total_movies
    similarity =  sum(mean_sim_scores_movies)/len(mean_sim_scores_movies
    diversity = 1- similarity

    return mean_sim_scores_movies, similarity, diversity, coverage,total_
```

# Findings for content-based recommendations

*Content based recommendatons - Similarity, Diversity, Coverage*

In [27]:
```python
content_similarity = compute_mean_sim_score(400,content_list_of_recommend
content_diversity = compute_mean_sim_score(400,content_list_of_recommenda
content_coverage = compute_mean_sim_score(400,content_list_of_recommendat
print("Collaborative similarity", content_similarity)
print("Collaborative  Diversity", content_diversity)
print("Collaborative  coverage", content_coverage)
```

```
Collaborative similarity 0.4648508515411847
Collaborative  Diversity 0.5351491484588153
Collaborative  coverage 0.5953333333333334
```

# Findings for Collaborative filtering

## Collaborative recommendatons - Similarity, Diversity, Coverage

In [28]:
```
1 collab_similarity = compute_mean_sim_score(400,collaborative_list_of_rec
2 collab_diversity = compute_mean_sim_score(400,collaborative_list_of_reco
3 collab_coverage = compute_mean_sim_score(400,collaborative_list_of_recom
4 print("Collaborative similarity", collab_similarity)
5 print("Collaborative  Diversity", collab_diversity)
6 print("Collaborative  coverage", collab_coverage)
7
8
```

```
Collaborative similarity 0.06625481199399376
Collaborative  Diversity 0.9337451880060063
Collaborative  coverage 0.6813333333333333
```
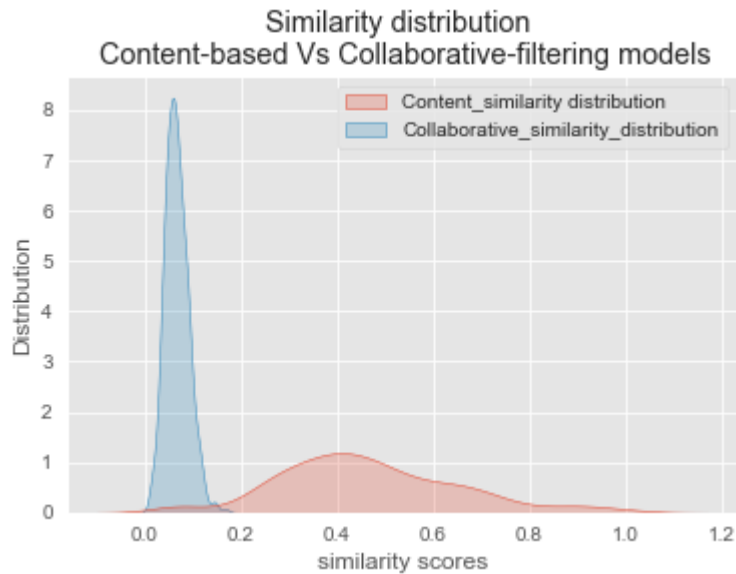
In [ ]:
```
1
```

## Similarity distribution comparison, content-based Vs Collaborative

In [29]:
```
1 content_mean_sim_scores= compute_mean_sim_score(400,content_list_of_reco
2 collab_mean_sim_scores = compute_mean_sim_score(400,collaborative_list_o
```

In [30]:
```
1 content_collab_sim_scores = pd.DataFrame({"Content_similarity distributi
2                          "Collaborative_similarity_distribution" : collab_mean_s
```

In [31]: ▶| 
```
1  sim_dis = sns.kdeplot(data = content_collab_sim_scores, fill =True, )
2  sim_dis.set_xlabel("similarity scores", fontsize = 12)
3  sim_dis.set_ylabel("Distribution", fontsize = 12)
4  sim_dis.set(title='Similarity distribution \n Content-based Vs Collaborat
```

Out[31]: [Text(0.5, 1.0, 'Similarity distribution \n Content-based Vs Collaborative-
filtering models')]

Similarity distribution
Content-based Vs Collaborative-filtering models

## Popularity plot comparison, content-based Vs Collaborative

In [34]: ▶|
```
1  from collections import Counter
2  total_content_recommendations = compute_mean_sim_score(400,content_list_c
3  content_based_count = [ ]
4  content_counter = Counter(total_content_recommendations)
5  for i in total_content_recommendations:
6      content_based_count.append(content_counter[i])
```
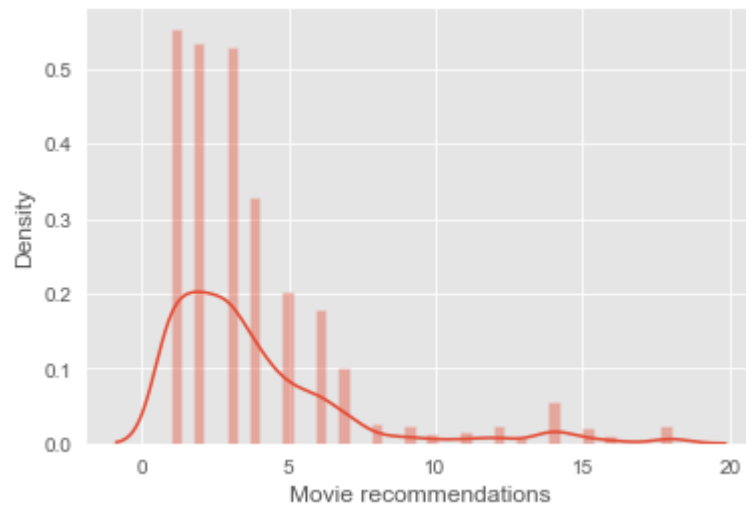
In [35]: ▶|
```
1  total_collab_recommendations = compute_mean_sim_score(400,collaborative_l
2  collab_based_count = []
3  collab_counter = Counter(total_collab_recommendations)
4  for i in total_collab_recommendations:
5      collab_based_count.append(collab_counter[i])
```

In [36]: ▶|
```
1  import warnings
2  warnings.filterwarnings('ignore')
```
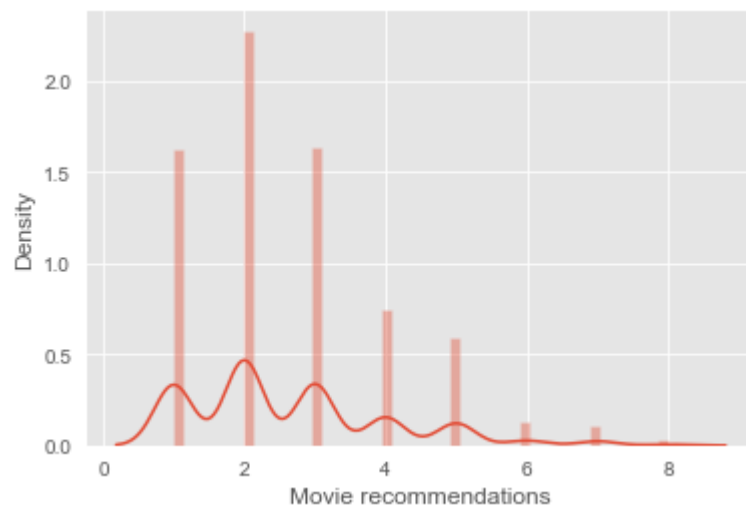
In [37]: ▶| 1 `sns.distplot(pd.Series(content_based_count, name = "Movie recommendations`

Out[37]: `<AxesSubplot:xlabel='Movie recommendations', ylabel='Density'>`



In [38]: ▶| 1 `sns.distplot(pd.Series(collab_based_count, name = "Movie recommendations`

Out[38]: `<AxesSubplot:xlabel='Movie recommendations', ylabel='Density'>`



------------------------------------------------------------------------------------
-------------------------------------------------

# Hybrid Recommendation

## Fetch highest rated movie for each user in the dataset to use as historical data for each user's content-based recommendation

```
In [39]:  1  highest_rated = ratings.groupby(['user_id'])['rating'].transform(max) ==
          2  highest_rated_df = ratings[highest_rated]
          3  total_highest_rated_movies = highest_rated_df.groupby(['user_id','movie_
          4  highest_rated_movies = total_highest_rated_movies.groupby('user_id').app
          5  highest_rated_movies['user_id'] = highest_rated_movies['user_id'].astype
          6  highest_rated_movies = highest_rated_movies.sort_values("user_id")
```

## Function to get user ID and the title of the movie they rated highly

```
In [40]:  1  def extract_users_movie(df, rating):
          2      dic ={}
          3      user = 1
          4      for movieid in rating['movie_id']:
          5          dic[user] = list(df.loc[df['movie_id'] == movieid, 'title'])[0]
          6          user = user +1
          7  #     dic = {k: list(v[1]) for k,v in dic.items()}
          8      return dic
          9
```

```
In [41]:  1  users_history = extract_users_movie(movies, highest_rated_movies)
```

## Modify content recommender, make new set of recommendation to take a dictionary

```
In [42]:  1  df.index = range(len(df))
          2  movies.index =range(len(df))
```

```
In [43]:  1  def modified_content_recommender(dic,n_items,df):
          2      recommendation_list = []
          3      for movie in dic.values():
          4          recommendation_list.append(list(content_recommender(df, str(movi
          5      return recommendation_list
          6
```

## Build a hybrid recommender

## New content-based recommender + the collaborative model from above

*Function will also return purely content and purely collaborative recommendaitons for comparison*

In [44]: ▶

```python
1  def hybrid_recommender(df,n_users,n_items, history, top_n_conent, top_n_
2      content_based = modified_content_recommender(history,n_items, df)
3      collaborative_filter = collaborative_recommender(n_users,df)
4      total_recommendations = []
5      for i in range(len(content_based )):
6          total_recommendations.append(list(set(random.sample(content_based
7      return total_recommendations, content_based, collaborative_filter
```

## Make all recommendations

In [46]: ▶

```python
1  result = hybrid_recommender(df,400,10,users_history,2,8,dpdf)
```

# Findings from all 3 models

*Content based recommendations analysis*

In [47]: ▶

```python
1  print("content based similarity", compute_mean_sim_score(400,result[1],d
2  print("content based Diversity",compute_mean_sim_score(400,result[1],dpd
3  print("content based coverage", compute_mean_sim_score(400,result[1],dpd
```

```
content based similarity 0.46492154055899637
content based Diversity 0.5350784594410036
content based coverage 0.581
```

*Collaborative filtering recommendations analysis*

In [48]: ▶

```python
1  print("collaborative similarity", compute_mean_sim_score(400,result[2],d
2  print("collaborative Diversity",compute_mean_sim_score(400,result[2],dpd
3  print("collaborative coverage", compute_mean_sim_score(400,result[2],dpd
```

```
collaborative similarity 0.06625481199399376
collaborative Diversity 0.9337451880060063
collaborative coverage 0.6813333333333333
```

*Hybrid based recommendations analysis*

```
In [49]: ▶| 1 print("Hybrid based similarity", compute_mean_sim_score(400,result[0],dp
          2 print("Hybrid based Diversity",compute_mean_sim_score(400,result[0],dpdf.
          3 print("Hybrid based coverage", compute_mean_sim_score(400,result[0],dpdf.
```

```
Hybrid based similarity 0.07296290370859333
Hybrid based Diversity 0.9270370962914066
Hybrid based coverage 0.6946666666666667
```

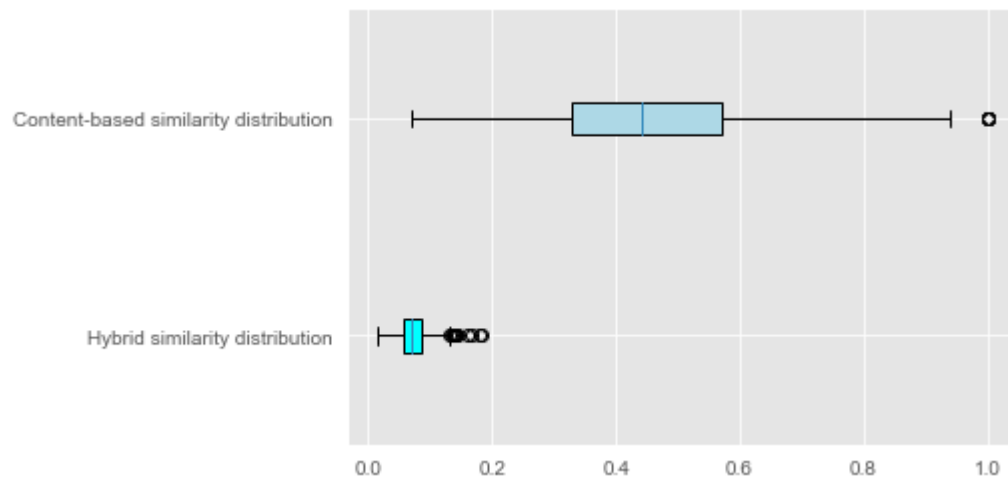## Similarity & diversity comparison Between Hybrid and content-based model

```
In [50]: ▶| 1 # compute_mean_sim_score(400,result[0],dpdf,movies)
          2 print("Mean sim score for Hybrid recommendations", (compute_mean_sim_sco
          3 print("Mean sim score for Content_based recommendations", (compute_mean_s
```

```
Mean sim score for Hybrid recommendations 0.07296290370859333
Mean sim score for Content_based recommendations 0.46492154055899637
```

```
In [51]: ▶| 1 # compute_mean_sim_score(400,result[0],dpdf,movies)
          2 print("Diversity for Hybrid recommendations", np.mean(compute_mean_sim_s
          3 print("Diversity for Content_based recommendations", np.mean(compute_mea
          4 print("% increase in diversity by Hybrid model",np.mean(compute_mean_sim
            ◀ ▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬                                                  ▶
```

```
Diversity for Hybrid recommendations 0.9270370962914066
Diversity for Content_based recommendations 0.5350784594410036
% increase in diversity by Hybrid model 0.391958636850403
```

```
1  box_plot_data=[compute_mean_sim_score(400,result[0],dpdf,movies)[0],comp
2  box= plt.boxplot(box_plot_data,vert=0,patch_artist=True,labels=['Hybrid
3            )
4  colors = ['cyan', 'lightblue']
5
6  for patch, color in zip(box['boxes'], colors):
7      patch.set_facecolor(color)
8
9  plt.show()
```



**Tests the significance in the difference of the two distributions using paired Student's t-Test**

```
1  from scipy.stats import ttest_rel
2  data1 = compute_mean_sim_score(400,result[0],dpdf, movies)[0]
3  data2 = compute_mean_sim_score(400,result[1],dpdf, movies)[0]
4  # compare samples
5  stat, p = ttest_rel(data1, data2)
6  print('Statistics=%.3f, p=%.20f' % (stat, p))
7  # interpret
8  alpha = 0.05
9  if p > alpha:
10     print('Same distributions (fail to reject H0)')
11 else:
12     print('Different distributions (reject H0)')
```

```
Statistics=-40.312, p=0.00000000000000000000
Different distributions (reject H0)
```

# Coverage comparison between Hybrid and Collaborative models

**% Difference in coverage, Hybrid vs Collaborative**

**At 0% data sparsity we there is 1.333 % difference in coverage b/n the hybrid and collaborative model**

In [54]: ▶| 1 `(compute_mean_sim_score(400,result[0],dpdf,movies)[3] - compute_mean_sim`

Out[54]: 1.3333333333333308

# Varying data sparsity

**Build new Hybrid and collaborative models to intake different Knn functions for eah step**

In [55]: ▶|
```python
1  def collaborative_recommender2(n_users,df, col_fun):
2      sample = df.copy()
3      list_of_recommendations=[]
4      for user in range(n_users+1):
5          if user == 0:
6              pass
7          else:
8              sample["estimated_rating"]= movies.apply(lambda x: col_fun.pr
9              list_of_recommendations.append(list(sample.sort_values('estir
10     return list_of_recommendations
11
```

In [56]: ▶|
```python
1  def hybrid_recommender2(df,n_users,n_items, history, top_n_conent, top_n
2      content_based = modified_content_recommender(history,n_items, df)
3      collaborative_filter = collaborative_recommender2(n_users,df,col_fun
4      total_recommendations = []
5      for i in range(len(content_based )):
6          total_recommendations.append(list(set(random.sample(content_base
7      return total_recommendations, content_based, collaborative_filter
```

*Varying data sparsity allows analysis of how the hybrid model performs compared to a purely collaborative model*

## For each step...

*1. Generate new sparse data for the selected level of sparsity*

*2. Train a new collaborative model*

*3. Input the new model into the hybrid model*

*4. Make purely content and collaborative recommendations*

*5. Make hybrid recommendations*

*6. Compare coverage between Hybrid and Collaborative model*

*7. Assess how diversity is impacted by the hybrid models in each steps*

## Coverage comparison at 16.7% sparsity

### Train the model with sparse data

In [57]:
```python
1  n_users = 401
2  n_movies = 2500
3  #generate a rating for each user/movie combination
4  data1 = pd.DataFrame(np.array(np.meshgrid([userid for userid in range(1,
5  np.random.seed(12)
6  randratings = np.random.randint(1,6, data1.shape[0])
7  data1['rating'] = randratings
8  data1["user_id"] = ratings["user_id"].astype(str)
9  data1["movie_id"] = ratings["movie_id"].astype(str)
10 data1.dtypes, data1.shape
```

Out[57]:
```
(user_id     object
 movie_id    object
 rating       int32
 dtype: object,
 (1000000, 3))
```

### Sparsity level

In [58]:
```python
1  (len(ratings) - len(data1))/len(ratings) # % sparcity
```

Out[58]:  0.16666666666666666

### Build new collaborative model on sparse data

```
In [59]:  ▶❙   1  our_seed = 14
              2  reader1 = Reader(rating_scale=(1,5)) # defaults to (0,5)
              3  d1 = Dataset.load_from_df(data1, reader1)
              4  knn1 = KNNBasic(k= 4, n_jobs=-1, verbose=False)
              5  random.seed(our_seed)
              6  np.random.seed(our_seed)
              7  knn_cv1 = cross_validate(knn1, d1, measures=['RMSE'], cv=5, verbose=True
              8  print(knn_cv1)
              9  knn_RMSE1 = np.mean(knn_cv1['test_rmse'])
             10  print(f'\nThe RMSE across five folds was {knn_RMSE1}')
```

```
Evaluating RMSE of algorithm KNNBasic on 5 split(s).

                  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5  Mean    Std
RMSE (testset)    1.5791  1.5809  1.5782  1.5802  1.5794  1.5796  0.0009
Fit time          12.39   10.72   11.71   10.03   12.30   11.43   0.92
Test time         52.35   41.82   41.44   40.28   47.06   44.59   4.53
{'test_rmse': array([1.57911029, 1.58091823, 1.57821139, 1.58022001, 1.5793
8965]), 'fit_time': (12.391668558120728, 10.724379777908325, 11.71349334716
7969, 10.031909227371216, 12.297916173934937), 'test_time': (52.34842324256
897, 41.821391105651855, 41.44197487831116, 40.28261470794678, 47.056489706
03943)}

The RMSE across five folds was 1.5795699139520485
```

### Extract highly rated movies for each user based on the new sparse data

```
In [60]:  ▶❙   1  highest_rated1 = data1.groupby(['user_id'])['rating'].transform(max) ==
              2  highest_rated_df1 = data1[highest_rated1]
              3  total_highest_rated_movies1 = highest_rated_df1.groupby(['user_id','movi
              4  highest_rated_movies1 = total_highest_rated_movies1.groupby('user_id').a
              5  highest_rated_movies1['user_id'] = highest_rated_movies1['user_id'].asty
              6  highest_rated_movies1 = highest_rated_movies1.sort_values("user_id")
```

```
In [61]:  ▶❙   1  users_history1 = extract_users_movie(movies, highest_rated_movies1)
```

### Make a hybrid recommendations based on new highly rated movies for content-based recommendations and newly trained collaborative model

```
In [62]:  ▶❙   1  result1 = hybrid_recommender2(df,400,10,users_history1,2,8,dpdf,knn1)
```

```
In [63]: ▶| 1  print("Collaborative coverage", compute_mean_sim_score(400,result1[2],dpo
         2  print("Hybrid based coverage", compute_mean_sim_score(400,result1[0],dpd-
         3  print("% difference in coverage", (compute_mean_sim_score(400,result1[0].
```

Collaborative coverage 0.62
Hybrid based coverage 0.6236666666666667
% difference in coverage 0.3666666666666707

## Coverage comparison at 33.3% sparsity

```
In [64]: ▶|  1  n_users = 401
             2  n_movies = 2000
             3  #generate a rating for each user/movie combination
             4  data2 = pd.DataFrame(np.array(np.meshgrid([userid for userid in range(1,
             5  np.random.seed(12)
             6  randratings = np.random.randint(1,6, data2.shape[0])
             7  data2['rating'] = randratings
             8  data2["user_id"] = data2["user_id"].astype(str)
             9  data2["movie_id"] = data2["movie_id"].astype(str)
            10  data2.dtypes, data2.shape
```

```
Out[64]: (user_id      object
          movie_id     object
          rating        int32
          dtype: object,
          (800000, 3))
```

```
In [65]: ▶|  1  (len(ratings)-len(data2))/len(ratings)
```

Out[65]: 0.3333333333333333

```
In [66]:  ▶|   1  our_seed = 14
              2  reader2 = Reader(rating_scale=(1,5)) # defaults to (0,5)
              3  d2 = Dataset.load_from_df(data2, reader2)
              4  knn2 = KNNBasic(k= 4, n_jobs=-1, verbose=False)
              5  random.seed(our_seed)
              6  np.random.seed(our_seed)
              7  knn_cv2 = cross_validate(knn2, d2, measures=['RMSE'], cv=5, verbose=True
              8  print(knn_cv2)
              9  knn_RMSE2 = np.mean(knn_cv2['test_rmse'])
             10  print(f'\nThe RMSE across five folds was {knn_RMSE2}')
```

```
Evaluating RMSE of algorithm KNNBasic on 5 split(s).

                  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5  Mean    Std
RMSE (testset)    1.5847  1.5816  1.5831  1.5819  1.5841  1.5831  0.0012
Fit time          11.47   12.07   12.52   13.00   11.88   12.19   0.53
Test time         42.78   40.74   39.97   41.89   40.48   41.17   1.02
{'test_rmse': array([1.58467458, 1.58160778, 1.58306427, 1.58192259, 1.5841
4689]), 'fit_time': (11.469969749450684, 12.067939281463623, 12.51714825630
188, 12.996087789535522, 11.882978200912476), 'test_time': (42.783443689346
31, 40.741952657699585, 39.972779273986816, 41.891682147979736, 40.47976732
254028)}

The RMSE across five folds was 1.5830832199358247
```

```
In [67]:  ▶|   1  highest_rated2 = data2.groupby(['user_id'])['rating'].transform(max) == (
              2  highest_rated_df2 = data2[highest_rated2]
              3  total_highest_rated_movies2 = highest_rated_df2.groupby(['user_id','movi
              4  highest_rated_movies2 = total_highest_rated_movies2.groupby('user_id').a
              5  highest_rated_movies2['user_id'] = highest_rated_movies2['user_id'].asty
              6  highest_rated_movies2 = highest_rated_movies2.sort_values("user_id")
```

```
In [68]:  ▶|   1  users_history2 = extract_users_movie(movies, highest_rated_movies2)
```

```
In [69]:  ▶|   1  result2 = hybrid_recommender2(df,400,10,users_history2,2,8,dpdf,knn2)
```

```
In [70]:  ▶|   1  print("Collaborative coverage", compute_mean_sim_score(400,result2[2],dp
              2  print("Hybrid based coverage", compute_mean_sim_score(400,result2[0],dpd
              3  print("% difference in coverage", (compute_mean_sim_score(400,result2[0]
```

```
Collaborative coverage 0.541
Hybrid based coverage 0.5986666666666667
% difference in coverage 5.766666666666664
```

## *Coverage comparison at 50 % sparsity*

```
In [71]:  ▶│   1  n_users = 401
              2  n_movies = 1500
              3  #generate a rating for each user/movie combination
              4  data3 = pd.DataFrame(np.array(np.meshgrid([userid for userid in range(1,
              5  np.random.seed(12)
              6  randratings = np.random.randint(1,6, data3.shape[0])
              7  data3['rating'] = randratings
              8  data3["user_id"] = data3["user_id"].astype(str)
              9  data3["movie_id"] = data3["movie_id"].astype(str)
             10  data3.dtypes, data3.shape
```

```
Out[71]:  (user_id      object
           movie_id     object
           rating        int32
           dtype: object,
           (600000, 3))
```

```
In [72]:  ▶│   1  (len(ratings) - len(data3))/len(ratings)
```

```
Out[72]:  0.5
```

```
In [73]:  ▶│   1  our_seed = 14
              2  reader3 = Reader(rating_scale=(1,5)) # defaults to (0,5)
              3  d3 = Dataset.load_from_df(data3, reader3)
              4  knn3 = KNNBasic(k= 4, n_jobs=-1, verbose=False) #the default for k is 40
              5  random.seed(our_seed)
              6  np.random.seed(our_seed)
              7  knn_cv3 = cross_validate(knn3, d3, measures=['RMSE'], cv=5, verbose=True
              8  print(knn_cv3)
              9  knn_RMSE3 = np.mean(knn_cv3['test_rmse'])
             10  print(f'\nThe RMSE across five folds was {knn_RMSE3}')
```

```
Evaluating RMSE of algorithm KNNBasic on 5 split(s).

                  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5  Mean    Std
RMSE (testset)    1.5761  1.5798  1.5808  1.5812  1.5827  1.5801  0.0022
Fit time          8.61    9.53    9.30    8.99    8.69    9.02    0.35
Test time         35.51   33.88   32.02   35.44   38.53   35.07   2.15
{'test_rmse': array([1.57611062, 1.57982727, 1.58080716, 1.58120509, 1.5826
6862]), 'fit_time': (8.612277030944824, 9.53093409538269, 9.29773187637329
1, 8.98525881767273, 8.688381433486938), 'test_time': (35.507365226745605,
33.877373456954956, 32.018059968948364, 35.43669366836548, 38.5349366664886
5)}

The RMSE across five folds was 1.580123751584545
```

```
In [74]:  ▶  1  highest_rated3 = data3.groupby(['user_id'])['rating'].transform(max) ==
              2  highest_rated_df3 = data3[highest_rated3]
              3  total_highest_rated_movies3 = highest_rated_df3.groupby(['user_id','movi
              4  highest_rated_movies3 = total_highest_rated_movies3.groupby('user_id').a
              5  highest_rated_movies3['user_id'] = highest_rated_movies3['user_id'].asty
              6  highest_rated_movies3 = highest_rated_movies3.sort_values("user_id")
```

```
In [75]:  ▶  1  users_history3 = extract_users_movie(movies, highest_rated_movies3)
```

```
In [76]:  ▶  1  result3 = hybrid_recommender2(df,400,10,users_history3,2,8,dpdf,knn3)
```

```
In [77]:  ▶  1  print("Collaborative coverage",compute_mean_sim_score(400,result3[2],dpd-
              2  print("Hybrid based coverage", compute_mean_sim_score(400,result3[0],dpd-
              3  print("% difference in coverage", (compute_mean_sim_score(400,result3[0].
```

```
Collaborative coverage 0.442
Hybrid based coverage 0.5393333333333333
% difference in coverage 9.733333333333333
```

## *Coverage comparison at 66.66 % sparsity*

```
In [78]:  ▶   1  n_users = 401
              2  n_movies = 1000
              3  #generate a rating for each user/movie combination
              4  data4 = pd.DataFrame(np.array(np.meshgrid([userid for userid in range(1,
              5  np.random.seed(12)
              6  randratings = np.random.randint(1,6, data4.shape[0])
              7  data4['rating'] = randratings
              8  data4["user_id"] = data4["user_id"].astype(str)
              9  data4["movie_id"] = data4["movie_id"].astype(str)
             10  data4.dtypes, data4.shape
             11  print("Sparsity level", (len(ratings) - len(data4))/len(ratings) * 100)
```

```
Sparsity level 66.66666666666666
```

```
In [79]:  ▶| 1  our_seed = 14
             2  reader4 = Reader(rating_scale=(1,5))
             3  d4 = Dataset.load_from_df(data4, reader4)
             4  knn4 = KNNBasic(k= 4, n_jobs=-1, verbose=False)
             5  random.seed(our_seed)
             6  np.random.seed(our_seed)
             7  knn_cv4 = cross_validate(knn4, d4, measures=['RMSE'], cv=5, verbose=True)
             8  print(knn_cv3)
             9  knn_RMSE4 = np.mean(knn_cv4['test_rmse'])
            10  print(f'\nThe RMSE across five folds was {knn_RMSE4}')
```

```
Evaluating RMSE of algorithm KNNBasic on 5 split(s).

                 Fold 1  Fold 2  Fold 3  Fold 4  Fold 5  Mean    Std
RMSE (testset)   1.5766  1.5767  1.5789  1.5841  1.5798  1.5792  0.0027
Fit time         6.31    6.29    6.19    6.07    6.20    6.21    0.09
Test time        22.06   23.89   21.13   22.10   22.89   22.42   0.92
{'test_rmse': array([1.57611062, 1.57982727, 1.58080716, 1.58120509, 1.5826
6862]), 'fit_time': (8.612277030944824, 9.53093409538269, 9.29773187637329
1, 8.98525881767273, 8.688381433486938), 'test_time': (35.507365226745605,
33.877373456954956, 32.018059968948364, 35.43669366836548, 38.5349366664886
5)}

The RMSE across five folds was 1.579211162818686
```

```
In [80]:  ▶| 1  highest_rated4 = data4.groupby(['user_id'])['rating'].transform(max) == (
             2  highest_rated_df4 = data4[highest_rated3]
             3  total_highest_rated_movies4 = highest_rated_df4.groupby(['user_id','movi
             4  highest_rated_movies4 = total_highest_rated_movies4.groupby('user_id').a
             5  highest_rated_movies4['user_id'] = highest_rated_movies4['user_id'].asty
             6  highest_rated_movies4 = highest_rated_movies4.sort_values("user_id")
```

```
In [81]:  ▶| 1  users_history4 = extract_users_movie(movies, highest_rated_movies4)
```

```
In [82]:  ▶| 1  result4 = hybrid_recommender2(df,400,10,users_history4,2,8,dpdf,knn4)
```

```
In [83]:  ▶| 1  print("Collaborative coverage",compute_mean_sim_score(400,result4[2],dpd
             2  print("Hybrid based coverage", compute_mean_sim_score(400,result4[0],dpd
             3  print("% difference in coverage", (compute_mean_sim_score(400,result4[0]
```

```
Collaborative coverage 0.31833333333333336
Hybrid based coverage 0.44633333333333336
% difference in coverage 12.8
```

# *Coverage comparison at 83.33 % sparsity*

```
In [84]:  1  n_users = 401
          2  n_movies = 500
          3  #generate a rating for each user/movie combination
          4  data5 = pd.DataFrame(np.array(np.meshgrid([userid for userid in range(1,
          5  np.random.seed(12)
          6  randratings = np.random.randint(1,6, data5.shape[0])
          7  data5['rating'] = randratings
          8  data5["user_id"] = data5["user_id"].astype(str)
          9  data5["movie_id"] = data5["movie_id"].astype(str)
         10  data5.dtypes, data4.shape
         11  print("Sparsity level", (len(ratings) - len(data5))/len(ratings) * 100)
```

Sparsity level 83.33333333333334

```
In [85]:  1  our_seed = 14
          2  reader5 = Reader(rating_scale=(1,5)) # defaults to (0,5)
          3  d5 = Dataset.load_from_df(data5, reader5)
          4  knn5 = KNNBasic(k= 4, n_jobs=-1, verbose=False)
          5  random.seed(our_seed)
          6  np.random.seed(our_seed)
          7  knn_cv5 = cross_validate(knn5, d5, measures=['RMSE'], cv=5, verbose=True
          8  print(knn_cv5)
          9  knn_RMSE5 = np.mean(knn_cv5['test_rmse'])
         10  print(f'\nThe RMSE across five folds was {knn_RMSE5}')
```

Evaluating RMSE of algorithm KNNBasic on 5 split(s).

```
                 Fold 1  Fold 2  Fold 3  Fold 4  Fold 5  Mean    Std
RMSE (testset)   1.5784  1.5812  1.5830  1.5794  1.5880  1.5820  0.0034
Fit time         2.56    2.92    2.70    2.88    2.81    2.77    0.13
Test time        10.14   10.56   11.18   10.07   10.07   10.41   0.43
{'test_rmse': array([1.5783887 , 1.58122628, 1.58296136, 1.57937938, 1.5879
8607]), 'fit_time': (2.5628745555877686, 2.9153695106506348, 2.696699619293
213, 2.8804097175598145, 2.8139145374298096), 'test_time': (10.141149759292
603, 10.559873580932617, 11.183897256851196, 10.071245908737183, 10.0743658
54263306)}
```

The RMSE across five folds was 1.5819883564771986

```
In [86]:  1  highest_rated5 = data5.groupby(['user_id'])['rating'].transform(max) ==
          2  highest_rated_df5 = data5[highest_rated5]
          3  total_highest_rated_movies5 = highest_rated_df5.groupby(['user_id','movie
          4  highest_rated_movies5 = total_highest_rated_movies5.groupby('user_id').a
          5  highest_rated_movies5['user_id'] = highest_rated_movies5['user_id'].asty
          6  highest_rated_movies5 = highest_rated_movies5.sort_values("user_id")
```

```
In [87]:  1  users_history5 = extract_users_movie(movies, highest_rated_movies5)
          2  result5 = hybrid_recommender2(df,400,10,users_history5,2,8,dpdf,knn5)
```

```
In [88]:  ▶|   1  print("Collaborative coverage",compute_mean_sim_score(400,result5[2],dpd·
              2  print("Hybrid based coverage", compute_mean_sim_score(400,result5[0],dpd·
              3  print("% difference in coverage", (compute_mean_sim_score(400,result5[0].
```

```
Collaborative coverage 0.16666666666666666
Hybrid based coverage 0.337
% difference in coverage 17.033333333333335
```

## *Coverage comparison at 91.66 % sparsity*

```
In [89]:  ▶|   1  n_users = 401
              2  n_movies = 250
              3  #generate a rating for each user/movie combination
              4  data6 = pd.DataFrame(np.array(np.meshgrid([userid for userid in range(1,
              5  np.random.seed(12)
              6  randratings = np.random.randint(1,6, data6.shape[0])
              7  data6['rating'] = randratings
              8  data6["user_id"] = data6["user_id"].astype(str)
              9  data6["movie_id"] = data6["movie_id"].astype(str)
             10  data6.dtypes, data6.shape
             11  print("Sparsity level", (len(ratings) - len(data6))/len(ratings) * 100)
```

```
Sparsity level 91.66666666666666
```

```
In [90]:  ▶|   1  our_seed = 14
              2  reader6 = Reader(rating_scale=(1,5))
              3  d6 = Dataset.load_from_df(data6, reader6)
              4  knn6 = KNNBasic(k= 4, n_jobs=-1, verbose=False)
              5  random.seed(our_seed)
              6  np.random.seed(our_seed)
              7  knn_cv6 = cross_validate(knn6, d6, measures=['RMSE'], cv=5, verbose=True
              8  print(knn_cv6)
              9  knn_RMSE6 = np.mean(knn_cv6['test_rmse'])
             10  print(f'\nThe RMSE across five folds was {knn_RMSE6}')
```

```
Evaluating RMSE of algorithm KNNBasic on 5 split(s).

                 Fold 1  Fold 2  Fold 3  Fold 4  Fold 5  Mean    Std
RMSE (testset)   1.5765  1.5722  1.5748  1.5883  1.5803  1.5784  0.0056
Fit time         1.32    1.24    1.23    1.39    1.20    1.28    0.07
Test time        4.49    4.66    5.31    4.59    4.43    4.70    0.32
{'test_rmse': array([1.57653474, 1.57224313, 1.57483776, 1.58830412, 1.5803
1657]), 'fit_time': (1.3204686641693115, 1.2416791915893555, 1.231664180755
6152, 1.392242193222046, 1.204784631729126), 'test_time': (4.48500418663024
9, 4.661531686782837, 5.309798002243042, 4.586774587631226, 4.4321067333221
436)}

The RMSE across five folds was 1.578447265269169
```

```
In [91]:  ▶|   1  highest_rated6 = data6.groupby(['user_id'])['rating'].transform(max) ==
              2  highest_rated_df6 = data6[highest_rated6]
              3  total_highest_rated_movies6 = highest_rated_df6.groupby(['user_id','movie
              4  highest_rated_movies6 = total_highest_rated_movies6.groupby('user_id').ap
              5  highest_rated_movies6['user_id'] = highest_rated_movies6['user_id'].asty
              6  highest_rated_movies6 = highest_rated_movies6.sort_values("user_id")
```

```
In [92]:  ▶|   1  users_history6 = extract_users_movie(movies, highest_rated_movies6)
              2  result6 = hybrid_recommender2(df,400,10,users_history6,2,8,dpdf,knn6)
```

```
In [93]:  ▶|   1  print("Collaborative coverage",compute_mean_sim_score(400,result6[2],dpd
              2  print("Hybrid based coverage", compute_mean_sim_score(400,result6[0],dpd
              3  print("% difference in coverage", (compute_mean_sim_score(400,result6[0]
```

```
Collaborative coverage 0.08333333333333333
Hybrid based coverage 0.2576666666666665
% difference in coverage 17.433333333333334
```

### *Collect coverage values*

```
In [94]:  ▶|   1  collaborative_coverage_values= [compute_mean_sim_score(400,result[2],dpd
              2                          compute_mean_sim_score(400,result1[2],dp
              3                          compute_mean_sim_score(400,result2[2],dp
              4                          compute_mean_sim_score(400,result3[2],dp
              5                          compute_mean_sim_score(400,result4[2],dp
              6                          compute_mean_sim_score(400,result5[2],dp
              7                          compute_mean_sim_score(400,result6[2],dp
              8
              9  Hybrid_coverage_values = [compute_mean_sim_score(400,result[0],dpdf, mov
             10                          compute_mean_sim_score(400,result1[0],dp
             11                          compute_mean_sim_score(400,result2[0],dp
             12                          compute_mean_sim_score(400,result3[0],dp
             13                          compute_mean_sim_score(400,result4[0],dp
             14                          compute_mean_sim_score(400,result5[0],dp
             15                          compute_mean_sim_score(400,result6[0],dp
             16  Coverage_difference = list()
             17  for item1, item2 in zip(Hybrid_coverage_values ,collaborative_coverage_va
             18      item = item1 - item2
             19      Coverage_difference.append(item* 100)
             20
             21  Sparsity_level = [0,16.7,33.3,50, 66.67, 83.33, 91.67]
             22  # Coverage_difference = Hybrid_coverage_values - collaborative_coverage_v
```
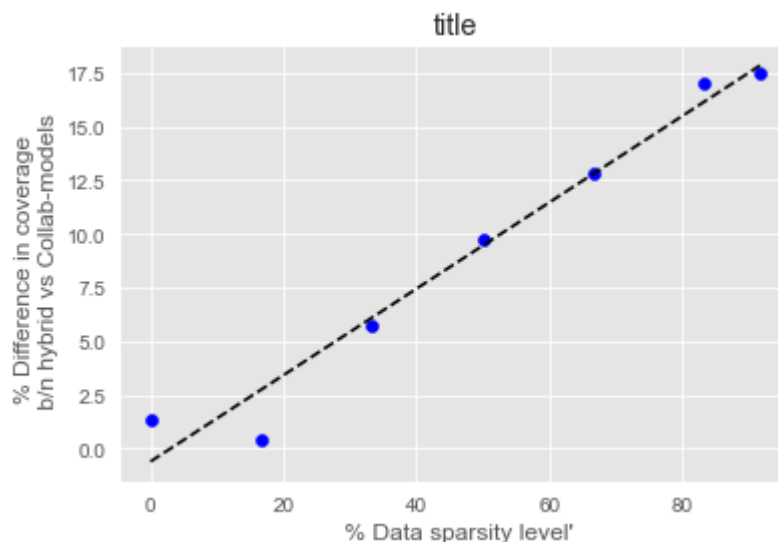
## Coverage comparison Hybrid Vs collaborative models for varying levels of data sparsity

```
In [125]:  ▶|   1  Coverage_data = pd.DataFrame(
               2      {'% Data sparsity level':Sparsity_level ,'Collaborative coverage': cc
               3       'Hybrid coverage': Hybrid_coverage_values ,
               4       '% Difference in coverage b/n hybrid model vs Collab-model': Coverage
               5      })
               6  Coverage_data
```

Out[125]:

|   | % Data sparsity level | Collaborative coverage | Hybrid coverage | % Difference in coverage b/n hybrid model vs Collab-model |
|---|---|---|---|---|
| 0 | 0.00 | 0.681333 | 0.694667 | 1.333333 |
| 1 | 16.70 | 0.620000 | 0.623667 | 0.366667 |
| 2 | 33.30 | 0.541000 | 0.598667 | 5.766667 |
| 3 | 50.00 | 0.442000 | 0.539333 | 9.733333 |
| 4 | 66.67 | 0.318333 | 0.446333 | 12.800000 |
| 5 | 83.33 | 0.166667 | 0.337000 | 17.033333 |
| 6 | 91.67 | 0.083333 | 0.257667 | 17.433333 |

```
In [96]:  ▶|   1  plt.scatter(Coverage_data['% Data sparsity level'], Coverage_data['% Dif
              2  z = np.polyfit(Coverage_data['% Data sparsity level'], Coverage_data['% [
              3  p = np.poly1d(z)
              4  plt.plot(Coverage_data['% Data sparsity level'],p(Coverage_data['% Data :
              5  plt.title("title")
              6  plt.xlabel(" % Data sparsity level'")
              7  plt.ylabel("% Difference in coverage \n b/n hybrid vs Collab-models")
              8  plt.show()
              9
             10  # plt.show()
```



**Diversity comparison Hybrid Vs content-based models for**

# varying levels of data sparsity

```
In [97]:  ▶| 1  content_diversity_values= [compute_mean_sim_score(400,result[1],dpdf, mov
          2                          compute_mean_sim_score(400,result1[1],dp
          3                          compute_mean_sim_score(400,result2[1],dp
          4                          compute_mean_sim_score(400,result3[1],dp
          5                          compute_mean_sim_score(400,result4[1],dp
          6                          compute_mean_sim_score(400,result5[1],dp
          7                          compute_mean_sim_score(400,result6[1],dp
          8
          9  Hybrid_diversity_values = [compute_mean_sim_score(400,result[0],dpdf, mov
         10                          compute_mean_sim_score(400,result1[0],dp
         11                          compute_mean_sim_score(400,result2[0],dp
         12                          compute_mean_sim_score(400,result3[0],dp
         13                          compute_mean_sim_score(400,result4[0],dp
         14                          compute_mean_sim_score(400,result5[0],dp
         15                          compute_mean_sim_score(400,result6[0],dp
         16  diversity_difference= list()
         17  for item1, item2 in zip(Hybrid_diversity_values ,content_diversity_value
         18      item = item1 - item2
         19      diversity_difference.append(item* 100)
         20
         21  Sparsity_level = [0,16.7,33.3,50, 66.67, 83.33, 91.67]
         22  # Coverage_difference = Hybrid_coverage_values - collaborative_coverage_v
```

```
In [98]:  ▶| 1  diversity_data = pd.DataFrame(
          2      {'% Data sparsity level':Sparsity_level ,'Content model diversity':
          3       'Hybrid model diversity': Hybrid_diversity_values,
          4       '% Difference in diversity': diversity_difference
          5      })
          6  diversity_data
```
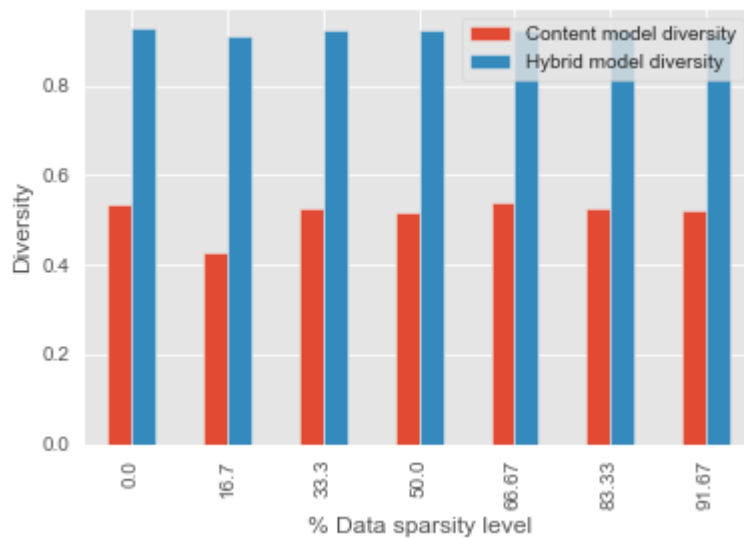
Out[98]:

|   | % Data sparsity level | Content model diversity | Hybrid model diversity | % Difference in diversity |
|---|---|---|---|---|
| **0** | 0.00 | 0.535078 | 0.927037 | 39.195864 |
| **1** | 16.70 | 0.428490 | 0.912906 | 48.441589 |
| **2** | 33.30 | 0.526261 | 0.925958 | 39.969780 |
| **3** | 50.00 | 0.518247 | 0.924863 | 40.661618 |
| **4** | 66.67 | 0.538587 | 0.922692 | 38.410487 |
| **5** | 83.33 | 0.526765 | 0.923430 | 39.666460 |
| **6** | 91.67 | 0.523503 | 0.917408 | 39.390465 |

In [99]: ▶| 1
2 diversity_data[['Content model diversity', 'Hybrid model diversity']].to_

Out[99]: array([[0.53507846, 0.9270371 ],
        [0.4284904 , 0.91290629],
        [0.52626051, 0.92595832],
        [0.5182469 , 0.92486308],
        [0.53858691, 0.92269178],
        [0.52676508, 0.92342969],
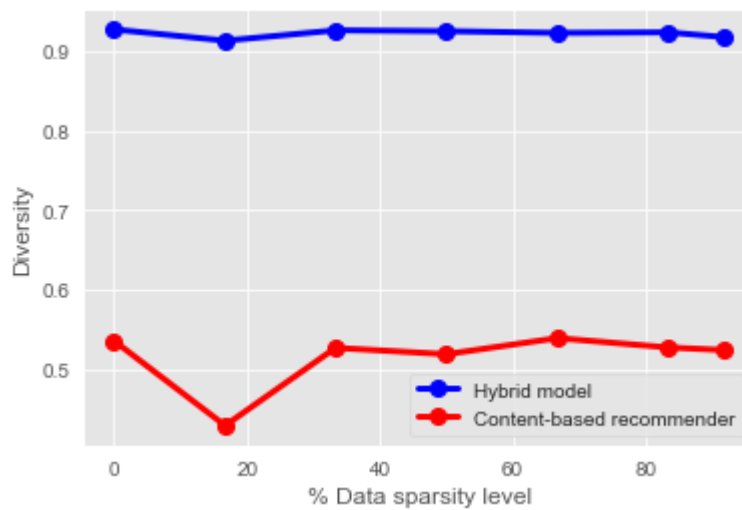        [0.5235031 , 0.91740776]])

In [100]: ▶| 1 diversity_data.plot(x= "% Data sparsity level", y =["Content model divers
          2

◀ ▓▓▓▓▓▓▓▓▓▓▓▓▓▓ ▶

Out[100]: Text(0, 0.5, 'Diversity')



**Better visual**

```
In [112]:  ▶  1  x = diversity_data['% Data sparsity level']
              2  y = diversity_data['Content model diversity']
              3  z = diversity_data['Hybrid model diversity']
              4
              5  # Plot a simple line chart
              6  plt.plot(x, z, 'b', label='Hybrid model', marker='o', markersize=8,linew:
              7
              8  plt.plot(x, y, 'r', label='Content-based recommender',marker='o',markers:
              9
             10  # Plot another line on the same chart/graph
             11  plt.xlabel('% Data sparsity level')
             12  plt.ylabel('Diversity')
             13  # plt.legend(loc='upper right')
             14  plt.legend()
             15  plt.show()
```



## Sample cosine similarity illustration

```
1  V = np.array([[1,5], [5,-10], [6,2], [-9,-4]])
2  origin = np.array([[0, 0, 0, 0],[0, 0, 0, 0]]) # origin point
3
4  plt.quiver(*origin, V[:,0], V[:,1], color=['r','b','g','black'], scale=3:
5  plt.ylabel('Romance')
6  plt.xlabel('Comedy')
7  plt.title('Romance or Comedy')
8  plt.show()
9
```