# Recommendations_with_IBM

December 21, 2020

# 1 Recommendations with IBM

In this notebook, you will be putting your recommendation skills to use on real data from the IBM Watson Studio platform.

You may either submit your notebook through the workspace here, or you may work from your local machine and submit through the next page. Either way assure that your code passes the project RUBRIC. **Please save regularly.**

By following the table of contents, you will build out a number of different methods for making recommendations that can be used for different situations.

## 1.1 Table of Contents

At the end of the notebook, you will find directions for how to submit your work. Let's get started by importing the necessary libraries and reading in the data.

```
In [56]: import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         import project_tests as t
         import pickle

         %matplotlib inline

         df = pd.read_csv('data/user-item-interactions.csv')
         df_content = pd.read_csv('data/articles_community.csv')
         del df['Unnamed: 0']
         del df_content['Unnamed: 0']

         # Show df to get an idea of the data
         df.head()

Out[56]:    article_id                                              title  \
         0      1430.0  using pixiedust for fast, flexible, and easier...
         1      1314.0         healthcare python streaming application demo
         2      1429.0            use deep learning for image classification
         3      1338.0             ml optimization using cognitive assistant
         4      1276.0                deploy your python model as a restful api
```

```
                                             email
       0   ef5f11f77ba020cd36e1105a00ab868bbdbf7fe7
       1   083cbdfa93c8444beaa4c5f5e0f5f9198e4f9e0b
       2   b96a4f2e92d8572034b1e9b28f9ac673765cd074
       3   06485706b34a5c9bf2a0ecdac41daf7e7654ceb7
       4   f01220c46fc92c6e6b161b1849de11faacd7ccb2
```

In [57]: df.shape

Out[57]: (45993, 3)

In [58]: df.describe()

```
Out[58]:           article_id
         count   45993.000000
         mean      908.846477
         std       486.647866
         min         0.000000
         25%       460.000000
         50%      1151.000000
         75%      1336.000000
         max      1444.000000
```

In [59]: df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 45993 entries, 0 to 45992
Data columns (total 3 columns):
article_id    45993 non-null float64
title         45993 non-null object
email         45976 non-null object
dtypes: float64(1), object(2)
memory usage: 1.1+ MB
```

In [60]: # Show df_content to get an idea of the data
         df_content.head()

```
Out[60]:                                                  doc_body  \
         0  Skip navigation Sign in SearchLoading...\r\n\r...
         1  No Free Hunch Navigation * kaggle.com\r\n\r\n ...
         2   * Login\r\n * Sign Up\r\n\r\n * Learning Pat...
         3  DATALAYER: HIGH THROUGHPUT, LOW LATENCY AT SCA...
         4  Skip navigation Sign in SearchLoading...\r\n\r...

                                             doc_description  \
         0  Detect bad readings in real time using Python ...
         1  See the forest, see the trees. Here lies the c...
```

```
         2   Heres this weeks news in Data Science and Bi...
         3   Learn how distributed DBs solve the problem of...
         4   This video demonstrates the power of IBM DataS...

                                         doc_full_name doc_status  article_id
         0   Detect Malfunctioning IoT Sensors with Streami...      Live           0
         1   Communicating data science: A guide to present...      Live           1
         2            This Week in Data Science (April 18, 2017)      Live           2
         3   DataLayer Conference: Boost the performance of...       Live           3
         4         Analyze NY Restaurant data using Spark in DSX      Live           4
```

In [61]: df_content.shape

Out[61]: (1056, 5)

In [62]: df.isnull().sum()

Out[62]: article_id     0
         title          0
         email         17
         dtype: int64

In [63]: df_content.describe()

Out[63]:         article_id
         count  1056.000000
         mean    523.913826
         std     303.480641
         min       0.000000
         25%     260.750000
         50%     523.500000
         75%     786.250000
         max    1050.000000

In [64]: df_content.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1056 entries, 0 to 1055
Data columns (total 5 columns):
doc_body          1042 non-null object
doc_description   1053 non-null object
doc_full_name     1056 non-null object
doc_status        1056 non-null object
article_id        1056 non-null int64
dtypes: int64(1), object(4)
memory usage: 41.3+ KB
```


In [65]: df_content.isnull().sum()

```
Out[65]: doc_body           14
         doc_description     3
         doc_full_name       0
         doc_status          0
         article_id          0
         dtype: int64
```

### 1.1.1  Part I : Exploratory Data Analysis

Use the dictionary and cells below to provide some insight into the descriptive statistics of the data.

1. What is the distribution of how many articles a user interacts with in the dataset? Provide a visual and descriptive statistics to assist with giving a look at the number of times each user interacts with an article.

```
In [66]: user_interaction = df.groupby('email').count()['article_id'].sort_values(ascending=Fals
         user_interaction

Out[66]: email
         2b6c0f514c2f2b04ad3c4583407dccd0810469ee    364
         77959baaa9895a7e2bdc9297f8b27c1b6f2cb52a    363
         2f5c7feae533ce046f2cb16fb3a29fe00528ed66    170
         a37adec71b667b297ed2440a9ff7dad427c7ac85    169
         8510a5010a5d4c89f5b07baac6de80cd12cfaf93    160
         f8c978bcf2ae2fb8885814a9b85ffef2f54c3c76    158
         284d0c17905de71e209b376e3309c0b08134f7e2    148
         d9032ff68d0fd45dfd18c0c5f7324619bb55362c    147
         18e7255ee311d4bd78f5993a9f09538e459e3fcc    147
         c60bb0a50c324dad0bffd8809d121246baef372b    145
         276d9d8ca0bf52c780b5a3fc554fa69e74f934a3    145
         56832a697cb6dbce14700fca18cffcced367057f    144
         b2d2c70ed5de62cf8a1d4ded7dd141cfbbdd0388    142
         ceef2a24a2a82031246814b73e029edba51e8ea9    140
         8dc8d7ec2356b1b106eb3d723f3c234e03ab3f1e    137
         e38f123afecb40272ba4c47cb25c96a9533006fa    136
         53db7ac77dbb80d6f5c32ed5d19c1a8720078814    132
         6c14453c049b1ef4737b08d56c480419794f91c2    131
         fd824fc62b4753107e3db7704cd9e8a4a1c961f1    116
         c45f9495a76bf95d2633444817f1be8205ad542d    114
         12bb8a9740400ced27ae5a7d4c990ac3b7e3c77d    104
         3427a5a4065625363e28ac8e85a57a9436010e9c    103
         497935037e41a94d2ae02488d098c7abda9a30bc    102
         0d644205ecefdef33e3346bb3551f5e68dc57c58    102
         e90de4b883d9de64a47774ad7ad49ca6fd69d4fe    101
         015aaf617598e413a35d6d2249e26b7f3c40adb7    101
         db1c400ffb74f14390deba2140bd31d2e1dc5c4e     98
         7dc02db8b76fffbdfe29542da672d4d5fd5ed4ae     97
         2e205a44014ca7bdbf07fc32f3c9d17699671d03     96
```

```
4070b8d82484ed99cdb9bbc2ebf4e9aca06fd934        95
                                               ...
42d4a9f766f2770e88a566cb65438a9b92446e6a         1
99a8fdeab6072b892f3477f2d91628df09cce12b         1
998ca3bffaaeb42f77cac8daf5f632a0c00b1c30         1
40002a2b20cee2d68bb9489ebd403ef9993100c2         1
9bbcd23976d1f9857fbb5e11291d37a2a2768341         1
9beb8742d40fb0619598cc3ae384165bca8d0794         1
efebe789cddce15baf08adab2c3da793896eb3cb         1
3e15c6b4972e54052ef3084190bdf1167b5db1a8         1
9db953fb65f5d57d8b8d82a0d04471dd5b7bac7b         1
9d3363969ba2a7f1d012d5c55af76652fc6ddc36         1
9d0375f208a9f91db408b5cf8da78e976fed3a55         1
9cfcf871ffb197ba5ad6bc6408ab5dc66d5b796d         1
9cfa28d68d71ba3fb1bf4745319be2258b87eb92         1
9ce6218339bd9186a3d0fe7da3494bc5af43dcba         1
9ce1e204a22ba4cd4a0a53da42238ae830b5879d         1
9cdb6449c080df01e366ce9c66f07a549be838d9         1
9cc6d232298678b4e24cf97ca0c74675fc2f132e         1
efe31a945040de5c0b5857b0072dc9254e96b37d         1
9c2394077e008013b92ec391eaf908d5ef3dd611         1
dc323e9b8ca2a9bf6397e43063fc093ae90788ea         1
9cb9845ca344b23b49ad94f4fddbcf95fedc0617         1
9cadbc14289d0db3937f00f4f2aab8d49b49680a         1
3f7be78857cda042074028beed41d088e5dd6a99         1
efded4d12cb4d1f53515e503d4ad3c4ca850a4da         1
3faaf951e4fa83cd67032688320d03d832ae708c         1
efdb4c363358224cd99d45053e2dbddf659e25ce         1
3fac88958dc7903b380743597f44a79cf76ea128         1
9c4b5dda1282c94128a7dc778951a313cce8055b         1
3fbe4978a20ee5ddc07648f2762b808ea18cedd1         1
6755c5d49a97e785583f65a92f72bc09459905a9         1
Name: article_id, Length: 5148, dtype: int64
```
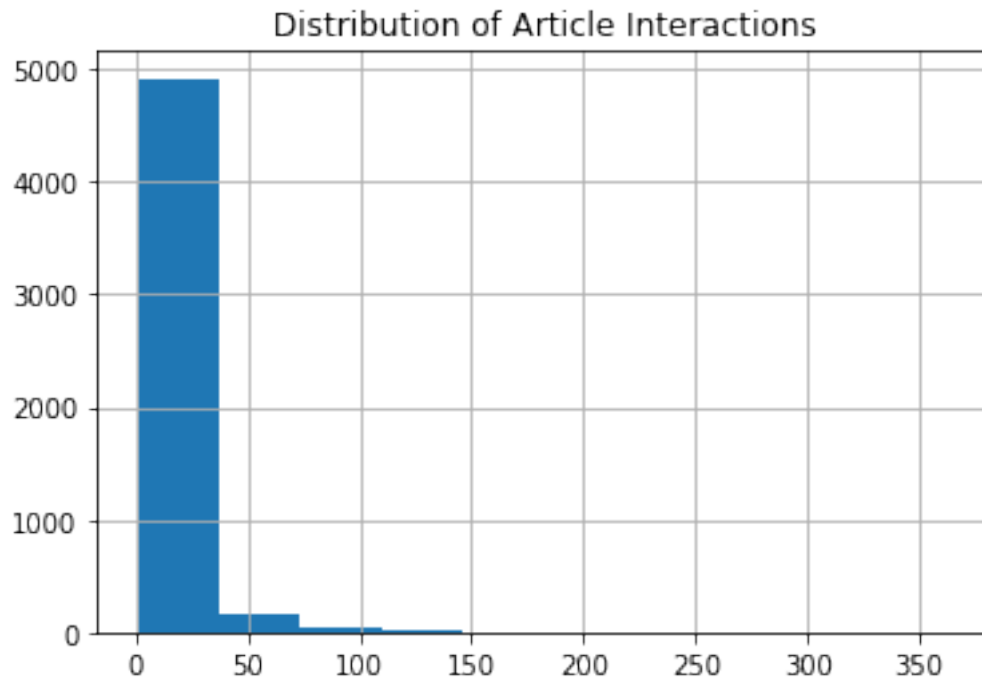
In [67]: # Visualization of User-articleS interaction
         user_interaction.hist()
         plt.title('Distribution of Article Interactions')

Out[67]: Text(0.5,1,'Distribution of Article Interactions')

Distribution of Article Interactions

In [68]: `# descriptive stats`
`user_interaction.describe()`

Out[68]: count     5148.000000
         mean         8.930847
         std         16.802267
         min          1.000000
         25%          1.000000
         50%          3.000000
         75%          9.000000
         max        364.000000
         Name: article_id, dtype: float64

In [69]: `user_interaction.median()`

Out[69]: 3.0

In [70]: `# Fill in the median and maximum number of user_article interactios below`

`median_val = 3 # 50% of individuals interact with ____ number of articles or fewer.`
`max_views_by_user = 364 # The maximum number of user-article interactions by any 1 user`

2. Explore and remove duplicate articles from the **df_content** dataframe.

In [71]: `# Find and explore duplicate articles`
`df_content.nunique()`

6

```
Out[71]: doc_body          1036
         doc_description   1022
         doc_full_name     1051
         doc_status           1
         article_id        1051
         dtype: int64
```

```
In [72]: df_content.duplicated("article_id").sum()
```

```
Out[72]: 5
```

```
In [73]: # Remove any rows that have the same article_id - only keep the first
         df_content.drop_duplicates(subset='article_id', keep='first', inplace=True)
```

```
In [74]: df_content.duplicated("article_id").sum()
```

```
Out[74]: 0
```

```
In [75]: df_content.shape
```

```
Out[75]: (1051, 5)
```

3. Use the cells below to find:
**a.** The number of unique articles that have an interaction with a user.
**b.** The number of unique articles in the dataset (whether they have any interactions or not). **c.** The number of unique users in the dataset. (excluding null values) **d.** The number of user-article interactions in the dataset.

```
In [76]: df.article_id.nunique()
```

```
Out[76]: 714
```

```
In [77]: df_content.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1051 entries, 0 to 1055
Data columns (total 5 columns):
doc_body          1037 non-null object
doc_description   1048 non-null object
doc_full_name     1051 non-null object
doc_status        1051 non-null object
article_id        1051 non-null int64
dtypes: int64(1), object(4)
memory usage: 49.3+ KB
```

```
In [78]: df.email.nunique()
```

```
Out[78]: 5148
```

7

```
In [79]: df.shape
```

```
Out[79]: (45993, 3)
```

```
In [80]: unique_articles = 714   # The number of unique articles that have at least one interacti
         total_articles = 1051  # The number of unique articles on the IBM platform
         unique_users = 5148 # The number of unique users
         user_article_interactions = 45993 # The number of user-article interactions
```

4. Use the cells below to find the most viewed **article_id**, as well as how often it was viewed. After talking to the company leaders, the email_mapper function was deemed a reasonable way to map users to ids. There were a small number of null values, and it was found that all of these null values likely belonged to a single user (which is how they are stored using the function below).

```
In [81]: df.groupby(["article_id"])["email"].count().sort_values(ascending=False).head()
```

```
Out[81]: article_id
         1429.0    937
         1330.0    927
         1431.0    671
         1427.0    643
         1364.0    627
         Name: email, dtype: int64
```

```
In [82]: most_viewed_article_id = "1429.0" # The most viewed article in the dataset as a string
         max_views = 937 # The most viewed article in the dataset was viewed how many times?
```

```
In [83]: df.head()
```

```
Out[83]:    article_id                                              title  \
         0     1430.0  using pixiedust for fast, flexible, and easier...
         1     1314.0         healthcare python streaming application demo
         2     1429.0           use deep learning for image classification
         3     1338.0            ml optimization using cognitive assistant
         4     1276.0               deploy your python model as a restful api

                                            email
         0  ef5f11f77ba020cd36e1105a00ab868bbdbf7fe7
         1  083cbdfa93c8444beaa4c5f5e0f5f9198e4f9e0b
         2  b96a4f2e92d8572034b1e9b28f9ac673765cd074
         3  06485706b34a5c9bf2a0ecdac41daf7e7654ceb7
         4  f01220c46fc92c6e6b161b1849de11faacd7ccb2
```

```
In [84]: ## No need to change the code here - this will be helpful for later parts of the notebo
         # Run this cell to map the user email to a user_id column and remove the email column

         def email_mapper():
             coded_dict = dict()
             cter = 1
```

```
        email_encoded = []

        for val in df['email']:
            if val not in coded_dict:
                coded_dict[val] = cter
                cter+=1

            email_encoded.append(coded_dict[val])
        return email_encoded

    email_encoded = email_mapper()
    del df['email']
    df['user_id'] = email_encoded

    # show header
    df.head()
```

```
Out[84]:    article_id                                         title  user_id
        0      1430.0  using pixiedust for fast, flexible, and easier...        1
        1      1314.0        healthcare python streaming application demo        2
        2      1429.0           use deep learning for image classification        3
        3      1338.0             ml optimization using cognitive assistant        4
        4      1276.0             deploy your python model as a restful api        5
```

```
In [85]: ## If you stored all your results in the variable names above,
         ## you shouldn't need to change anything in this cell

         sol_1_dict = {
             '`50% of individuals have _____ or fewer interactions.`': median_val,
             '`The total number of user-article interactions in the dataset is _____.`': user_a
             '`The maximum number of user-article interactions by any 1 user is _____.`': max_v
             '`The most viewed article in the dataset was viewed _____ times.`': max_views,
             '`The article_id of the most viewed article is _____.`': most_viewed_article_id,
             '`The number of unique articles that have at least 1 rating _____.`': unique_artic
             '`The number of unique users in the dataset is _____`': unique_users,
             '`The number of unique articles on the IBM platform`': total_articles
         }

         # Test your dictionary against the solution
         t.sol_1_test(sol_1_dict)
```

It looks like you have everything right here! Nice job!


### 1.1.2 Part II: Rank-Based Recommendations

Unlike in the earlier lessons, we don't actually have ratings for whether a user liked an article or not. We only know that a user has interacted with an article. In these cases, the popularity of an article can really only be based on how often an article was interacted with.

9

1. Fill in the function below to return the **n** top articles ordered with most interactions as the top. Test your function using the tests below.

```
In [86]: def get_top_articles(n, df=df):
             '''
             INPUT:
             n - (int) the number of top articles to return
             df - (pandas dataframe) df as defined at the top of the notebook

             OUTPUT:
             top_articles - (list) A list of the top 'n' article titles

             '''
             # Your code here

             top_articles = list(df.groupby(['title'])['article_id'].count().sort_values(ascendi

             return top_articles # Return the top article titles from df (not df_content)

         def get_top_article_ids(n, df=df):
             '''
             INPUT:
             n - (int) the number of top articles to return
             df - (pandas dataframe) df as defined at the top of the notebook

             OUTPUT:
             top_articles - (list) A list of the top 'n' article titles

             '''
             # Your code here
             top_articles = list(df['article_id'].value_counts().head(n).index)

             return top_articles # Return the top article ids

In [87]: print(get_top_articles(10))
         print(get_top_article_ids(10))

['use deep learning for image classification', 'insights from new york car accident reports', 'v
[1429.0, 1330.0, 1431.0, 1427.0, 1364.0, 1314.0, 1293.0, 1170.0, 1162.0, 1304.0]


In [88]: # Test your function by returning the top 5, 10, and 20 articles
         top_5 = get_top_articles(5)
         top_10 = get_top_articles(10)
         top_20 = get_top_articles(20)

         # Test each of your three lists from above
         t.sol_2_test(get_top_articles)
```

10

```
Your top_5 looks like the solution list! Nice job.
Your top_10 looks like the solution list! Nice job.
Your top_20 looks like the solution list! Nice job.
```

### 1.1.3 Part III: User-User Based Collaborative Filtering

1. Use the function below to reformat the **df** dataframe to be shaped with users as the rows and articles as the columns.

- Each **user** should only appear in each **row** once.

- Each **article** should only show up in one **column**.

- **If a user has interacted with an article, then place a 1 where the user-row meets for that article-column**. It does not matter how many times a user has interacted with the article, all entries where a user has interacted with an article should be a 1.

- **If a user has not interacted with an item, then place a zero where the user-row meets for that article-column**.

Use the tests to make sure the basic structure of your matrix matches what is expected by the solution.

```
In [89]: # create the user-article matrix with 1's and 0's

         def create_user_item_matrix(df):
             '''
             INPUT:
             df - pandas dataframe with article_id, title, user_id columns

             OUTPUT:
             user_item - user item matrix

             Description:
             Return a matrix with user ids as rows and article ids on the columns with 1 values
             an article and a 0 otherwise
             '''
             # Fill in the function here
             user_item=df.groupby(by=['user_id', 'article_id']).agg(lambda x: 1).unstack().filln


             return user_item # return the user_item matrix

         user_item = create_user_item_matrix(df)

In [90]: ## Tests: You should just need to run this cell.  Don't change the code.
         assert user_item.shape[0] == 5149, "Oops!  The number of users in the user-article matr
         assert user_item.shape[1] == 714, "Oops!  The number of articles in the user-article ma
         assert user_item.sum(axis=1)[1] == 36, "Oops!  The number of articles seen by user 1 do
         print("You have passed our quick tests!  Please proceed!")
```

11

You have passed our quick tests!  Please proceed!


In [91]: user_item.head()

Out[91]:              title                                                    \
        article_id 0.0    2.0    4.0    8.0    9.0    12.0   14.0   15.0   16.0
        user_id
        1               0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0
        2               0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0
        3               0.0    0.0    0.0    0.0    0.0    1.0    0.0    0.0    0.0
        4               0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0
        5               0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0


                        ...                                                    \
        article_id 18.0      ...    1434.0 1435.0 1436.0 1437.0 1439.0 1440.0 1441.0
        user_id              ...
        1               0.0  ...       0.0    0.0    1.0    0.0    1.0    0.0    0.0
        2               0.0  ...       0.0    0.0    0.0    0.0    0.0    0.0    0.0
        3               0.0  ...       0.0    0.0    1.0    0.0    0.0    0.0    0.0
        4               0.0  ...       0.0    0.0    0.0    0.0    0.0    0.0    0.0
        5               0.0  ...       0.0    0.0    0.0    0.0    0.0    0.0    0.0


        article_id 1442.0 1443.0 1444.0
        user_id
        1               0.0    0.0    0.0
        2               0.0    0.0    0.0
        3               0.0    0.0    0.0
        4               0.0    0.0    0.0
        5               0.0    0.0    0.0

        [5 rows x 714 columns]

2. Complete the function below which should take a user_id and provide an ordered list of the most similar users to that user (from most similar to least similar). The returned result should not contain the provided user_id, as we know that each user is similar to him/herself. Because the results for each user here are binary, it (perhaps) makes sense to compute similarity as the dot product of two users.

Use the tests to test your function.

```
In [97]: def find_similar_users(user_id, user_item=user_item):
             '''
             INPUT:
             user_id - (int) a user_id
             user_item - (pandas dataframe) matrix of users by articles:
                         1's when a user has interacted with an article, 0 otherwise

             OUTPUT:
```

```
            similar_users - (list) an ordered list where the closest users (largest dot product
                            are listed first

            Description:
            Computes the similarity of every pair of users based on the dot product
            Returns an ordered

            '''
            # compute similarity of each user to the provided user
            similarity = {}
            for uid in user_item.index:
                similarity[uid] = np.dot(user_item.loc[user_id, :], user_item.loc[uid, :])



            # sort by similarity
            similarity_sort = sorted(similarity.items(), key=lambda kv: kv[1], reverse=True)



            # create list of just the ids
            most_similar_users = [key for (key, value) in similarity_sort]



            # remove the own user's id
            most_similar_users.remove(user_id)

            return most_similar_users # return a list of the users in order from most to least

In [98]: # Do a spot check of your function
         print("The 10 most similar users to user 1 are: {}".format(find_similar_users(1)[:10]))
         print("The 5 most similar users to user 3933 are: {}".format(find_similar_users(3933)[:
         print("The 3 most similar users to user 46 are: {}".format(find_similar_users(46)[:3]))

The 10 most similar users to user 1 are: [3933, 23, 3782, 203, 4459, 131, 3870, 46, 4201, 49]
The 5 most similar users to user 3933 are: [1, 23, 3782, 203, 4459]
The 3 most similar users to user 46 are: [4201, 23, 3782]
```

3. Now that you have a function that provides the most similar users to each user, you will want to use these users to find articles you can recommend. Complete the functions below to return the articles you would recommend to each user.

```
In [99]: def get_article_names(article_ids, df=df):
             '''
             INPUT:
             article_ids - (list) a list of article ids
             df - (pandas dataframe) df as defined at the top of the notebook

             OUTPUT:
             article_names - (list) a list of article names associated with the list of article
```

```python
                            (this is identified by the title column)
    '''
    # Your code here
    #article_names = df[df['article_id'].isin(article_ids)]['title'].unique().tolist()
    article_names = df[df['article_id'].isin(article_ids)]['title'].drop_duplicates().v



    return article_names # Return the article names associated with list of article ids


def get_user_articles(user_id, user_item=user_item):
    '''
    INPUT:
    user_id - (int) a user id
    user_item - (pandas dataframe) matrix of users by articles:
                1's when a user has interacted with an article, 0 otherwise

    OUTPUT:
    article_ids - (list) a list of the article ids seen by the user
    article_names - (list) a list of article names associated with the list of article
                    (this is identified by the doc_full_name column in df_content)

    Description:
    Provides a list of the article_ids and article titles that have been seen by a user
    '''
    # Your code here

    article_ids = [str(id) for id in list(user_item.loc[user_id][user_item.loc[user_id]
    article_names = get_article_names(article_ids)



    return article_ids, article_names # return the ids and names


def user_user_recs(user_id, m=10):
    '''
    INPUT:
    user_id - (int) a user id
    m - (int) the number of recommendations you want for the user

    OUTPUT:
    recs - (list) a list of recommendations for the user

    Description:
    Loops through the users based on closeness to the input user_id
    For each user - finds articles the user hasn't seen before and provides them as rec
    Does this until m recommendations are found
```

```
        Notes:
        Users who are the same closeness are chosen arbitrarily as the 'next' user

        For the user where the number of recommended articles starts below m
        and ends exceeding m, the last items are chosen arbitrarily

        '''
        # Your code here
        recs = []
        most_similar_users = find_similar_users(user_id)
        viewed_article_ids_self, viewed_article_names_self = get_user_articles(user_id)


        for user in most_similar_users:
            article_ids, article_names = get_user_articles(user)

            for article_id in article_ids:
                if article_id not in viewed_article_ids_self:
                    if article_id not in recs and len(recs) < m:
                        recs.append(article_id)
                        if len(recs) >= m:
                            break
                if len(recs) >= m:
                    break


        return recs # return your recommendations for this user_id
```

In [100]: `# Check Results`
`get_article_names(user_user_recs(1, 10)) # Return 10 recommendations for user 1`

Out[100]: ['got zip code data? prep it for analytics.  ibm watson data lab  medium',
'timeseries data analysis of iot events by using jupyter notebook',
'graph-based machine learning',
'using brunel in ipython/jupyter notebooks',
'experience iot with coursera',
'the 3 kinds of context: machine learning and the art of the frame',
'deep forest: towards an alternative to deep neural networks',
'this week in data science (april 18, 2017)',
'higher-order logistic regression for large datasets',
'using machine learning to predict parking difficulty']

In [101]: `# Test your functions here - No need to change this code - just run this cell`
`assert set(get_article_names(['1024.0', '1176.0', '1305.0', '1314.0', '1422.0', '1427.`
`assert set(get_article_names(['1320.0', '232.0', '844.0'])) == set(['housing (2015): u`
`assert set(get_user_articles(20)[0]) == set(['1320.0', '232.0', '844.0'])`
`assert set(get_user_articles(20)[1]) == set(['housing (2015): united states demographi`

```
        assert set(get_user_articles(2)[0]) == set(['1024.0', '1176.0', '1305.0', '1314.0', '1
        assert set(get_user_articles(2)[1]) == set(['using deep learning to reconstruct high-r
        print("If this is all you see, you passed all of our tests!  Nice job!")
```

If this is all you see, you passed all of our tests!  Nice job!

4. Now we are going to improve the consistency of the **user_user_recs** function from above.

- Instead of arbitrarily choosing when we obtain users who are all the same closeness to a given user - choose the users that have the most total article interactions before choosing those with fewer article interactions.

- Instead of arbitrarily choosing articles from the user where the number of recommended articles starts below m and ends exceeding m, choose articles with the articles with the most total interactions before choosing those with fewer total interactions. This ranking should be what would be obtained from the **top_articles** function you wrote earlier.

```python
In [102]: def get_top_sorted_users(user_id, df=df, user_item=user_item):
              '''
              INPUT:
              user_id - (int)
              df - (pandas dataframe) df as defined at the top of the notebook
              user_item - (pandas dataframe) matrix of users by articles:
                      1's when a user has interacted with an article, 0 otherwise


              OUTPUT:
              neighbors_df - (pandas dataframe) a dataframe with:
                            neighbor_id - is a neighbor user_id
                            similarity - measure of the similarity of each user to the provide
                            num_interactions - the number of articles viewed by the user - if

              Other Details - sort the neighbors_df by the similarity and then by number of inte
                            highest of each is higher in the dataframe

              '''
              # Your code here
              neighbors_df = pd.DataFrame(columns=['neighbor_id', 'similarity', 'num_interaction
              for user in user_item.index:
                  if user == user_id:
                      continue
                  neighbors_df.loc[user] = [user, np.dot(user_item.loc[user_id, :], user_item.lo
                                          df[df['user_id']==user]['article_id'].count()]

              neighbors_df.sort_values(by=['similarity', 'num_interactions'], ascending=False, i

              return neighbors_df # Return the dataframe specified in the doc_string
```

16

```python
def user_user_recs_part2(user_id, m=10):
    '''
    INPUT:
    user_id - (int) a user id
    m - (int) the number of recommendations you want for the user

    OUTPUT:
    recs - (list) a list of recommendations for the user by article id
    rec_names - (list) a list of recommendations for the user by article title

    Description:
    Loops through the users based on closeness to the input user_id
    For each user - finds articles the user hasn't seen before and provides them as re
    Does this until m recommendations are found

    Notes:
    * Choose the users that have the most total article interactions
    before choosing those with fewer article interactions.

    * Choose articles with the articles with the most total interactions
    before choosing those with fewer total interactions.

    '''
    # Your code here
    recs = []

    neighbors_df = get_top_sorted_users(user_id)

    the_user_articles, the_article_names = get_user_articles(user_id)
    for user in neighbors_df['neighbor_id']:
        article_ids, article_names = get_user_articles(user)
        for id in article_ids:
            if id not in the_user_articles:
                recs.append(id)
            if len(recs) >= m:
                break
        if len(recs) >= m:
            break

    if len(recs) < m:
        for id in [str(id) for id in get_top_article_ids(100)]:
            if id not in the_user_articles:
                recs.append(id)
            if len(recs) >= m:
                break

    rec_names = get_article_names(recs)
```

```
                return recs, rec_names

In [103]:  # Quick spot check - don't change this code - just use it to test your functions
           rec_ids, rec_names = user_user_recs_part2(20, 10)
           print("The top 10 recommendations for user 20 are the following article ids:")
           print(rec_ids)
           print()
           print("The top 10 recommendations for user 20 are the following article names:")
           print(rec_names)

The top 10 recommendations for user 20 are the following article ids:
['12.0', '109.0', '125.0', '142.0', '164.0', '205.0', '302.0', '336.0', '362.0', '465.0']

The top 10 recommendations for user 20 are the following article names:
['timeseries data analysis of iot events by using jupyter notebook', 'dsx: hybrid mode', 'accele
```

5. Use your functions from above to correctly fill in the solutions to the dictionary below. Then test your dictionary against the solution. Provide the code you need to answer each following the comments below.

```
In [105]:  get_top_sorted_users(1).iloc[0]

Out[105]:  neighbor_id          3933.0
           similarity             35.0
           num_interactions       45.0
           Name: 3933, dtype: float64

In [107]:  get_top_sorted_users(1).neighbor_id.values[0]

Out[107]:  3933.0

In [106]:  get_top_sorted_users(131).iloc[9]

Out[106]:  neighbor_id           242.0
           similarity             25.0
           num_interactions      148.0
           Name: 242, dtype: float64

In [110]:  ### Tests with a dictionary of results

           user1_most_sim = 3933 # Find the user that is most similar to user 1
           user131_10th_sim = 242 # Find the 10th most similar user to user 131

In [111]:  ## Dictionary Test Here
           sol_5_dict = {
               'The user that is most similar to user 1.': user1_most_sim,
               'The user that is the 10th most similar to user 131': user131_10th_sim,
           }

           t.sol_5_test(sol_5_dict)
```

```
This all looks good!  Nice job!
```

6. If we were given a new user, which of the above functions would you be able to use to make recommendations? Explain. Can you think of a better way we might make recommendations? Use the cell below to explain a better method for new users.

**Provide your response here.**

7. Using your existing functions, provide the top 10 recommended articles you would provide for the a new user below. You can test your function against our thoughts to make sure we are all on the same page with how we might make a recommendation.

```
In [114]: new_user = '0.0'

          # What would your recommendations be for this new user '0.0'?  As a new user, they hav
          # Provide a list of the top 10 article ids you would give to
          new_user_recs = get_top_article_ids(10) # Your recommendations here
          new_user_recs = [str(ids) for ids in get_top_article_ids(10)]

In [115]: assert set(new_user_recs) == set(['1314.0','1429.0','1293.0','1427.0','1162.0','1364.0

          print("That's right!  Nice job!")
```

```
That's right!  Nice job!
```

### 1.1.4 Part IV: Content Based Recommendations (EXTRA - NOT REQUIRED)

Another method we might use to make recommendations is to perform a ranking of the highest ranked articles associated with some term. You might consider content to be the **doc_body**, **doc_description**, or **doc_full_name**. There isn't one way to create a content based recommendation, especially considering that each of these columns hold content related information.

1. Use the function body below to create a content based recommender. Since there isn't one right answer for this recommendation tactic, no test functions are provided. Feel free to change the function inputs if you decide you want to try a method that requires more input values. The input values are currently set with one idea in mind that you may use to make content based recommendations. One additional idea is that you might want to choose the most popular recommendations that meet your 'content criteria', but again, there is a lot of flexibility in how you might make these recommendations.

### 1.1.5 This part is NOT REQUIRED to pass this project. However, you may choose to take this on as an extra way to show off your skills.

```
In [ ]: def make_content_recs():
            '''
            INPUT:

            OUTPUT:

            '''
```

2. Now that you have put together your content-based recommendation system, use the cell below to write a summary explaining how your content based recommender works. Do you see any possible improvements that could be made to your function? Is there anything novel about your content based recommender?

### 1.1.6 This part is NOT REQUIRED to pass this project. However, you may choose to take this on as an extra way to show off your skills.

**Write an explanation of your content based recommendation system here.**

3. Use your content-recommendation system to make recommendations for the below scenarios based on the comments. Again no tests are provided here, because there isn't one right answer that could be used to find these content based recommendations.

### 1.1.7 This part is NOT REQUIRED to pass this project. However, you may choose to take this on as an extra way to show off your skills.

```
In [ ]: # make recommendations for a brand new user



        # make a recommendations for a user who only has interacted with article id '1427.0'
```

### 1.1.8 Part V: Matrix Factorization

In this part of the notebook, you will build use matrix factorization to make article recommendations to the users on the IBM Watson Studio platform.

1. You should have already created a **user_item** matrix above in **question 1** of **Part III** above. This first question here will just require that you run the cells to get things set up for the rest of **Part V** of the notebook.

```
In [116]: # Load the matrix here
          user_item_matrix = pd.read_pickle('user_item_matrix.p')

In [117]: # quick look at the matrix
          user_item_matrix.head()
```

```
Out[117]: article_id  0.0  100.0  1000.0  1004.0  1006.0  1008.0  101.0  1014.0  1015.0  \
          user_id
          1           0.0   0.0     0.0     0.0     0.0     0.0    0.0     0.0     0.0
          2           0.0   0.0     0.0     0.0     0.0     0.0    0.0     0.0     0.0
          3           0.0   0.0     0.0     0.0     0.0     0.0    0.0     0.0     0.0
          4           0.0   0.0     0.0     0.0     0.0     0.0    0.0     0.0     0.0
          5           0.0   0.0     0.0     0.0     0.0     0.0    0.0     0.0     0.0

          article_id  1016.0  ...   977.0  98.0  981.0  984.0  985.0  986.0  990.0  \
          user_id             ...
          1            0.0    ...    0.0   0.0    1.0    0.0    0.0    0.0    0.0
          2            0.0    ...    0.0   0.0    0.0    0.0    0.0    0.0    0.0
          3            0.0    ...    1.0   0.0    0.0    0.0    0.0    0.0    0.0
          4            0.0    ...    0.0   0.0    0.0    0.0    0.0    0.0    0.0
```

```
5                0.0  ...      0.0    0.0    0.0    0.0    0.0    0.0    0.0

article_id  993.0  996.0  997.0
user_id
1                0.0    0.0    0.0
2                0.0    0.0    0.0
3                0.0    0.0    0.0
4                0.0    0.0    0.0
5                0.0    0.0    0.0

[5 rows x 714 columns]
```

2. In this situation, you can use Singular Value Decomposition from numpy on the user-item matrix. Use the cell to perform SVD, and explain why this is different than in the lesson.

In [118]: *# Perform SVD on the User-Item Matrix Here*

```
u, s, vt = np.linalg.svd(user_item_matrix) # use the built in to get the three matrice
```

**Provide your response here.**

3. Now for the tricky part, how do we choose the number of latent features to use? Running the below cell, you can see that as the number of latent features increases, we obtain a lower error rate on making predictions for the 1 and 0 values in the user-item matrix. Run the cell below to get an idea of how the accuracy improves as we increase the number of latent features.

In [119]:
```python
num_latent_feats = np.arange(10,700+10,20)
sum_errs = []

for k in num_latent_feats:
    # restructure with k latent features
    s_new, u_new, vt_new = np.diag(s[:k]), u[:, :k], vt[:k, :]

    # take dot product
    user_item_est = np.around(np.dot(np.dot(u_new, s_new), vt_new))

    # compute error for each prediction to actual value
    diffs = np.subtract(user_item_matrix, user_item_est)

    # total errors and keep track of them
    err = np.sum(np.sum(np.abs(diffs)))
    sum_errs.append(err)


plt.plot(num_latent_feats, 1 - np.array(sum_errs)/df.shape[0]);
plt.xlabel('Number of Latent Features');
plt.ylabel('Accuracy');
plt.title('Accuracy vs. Number of Latent Features');
```
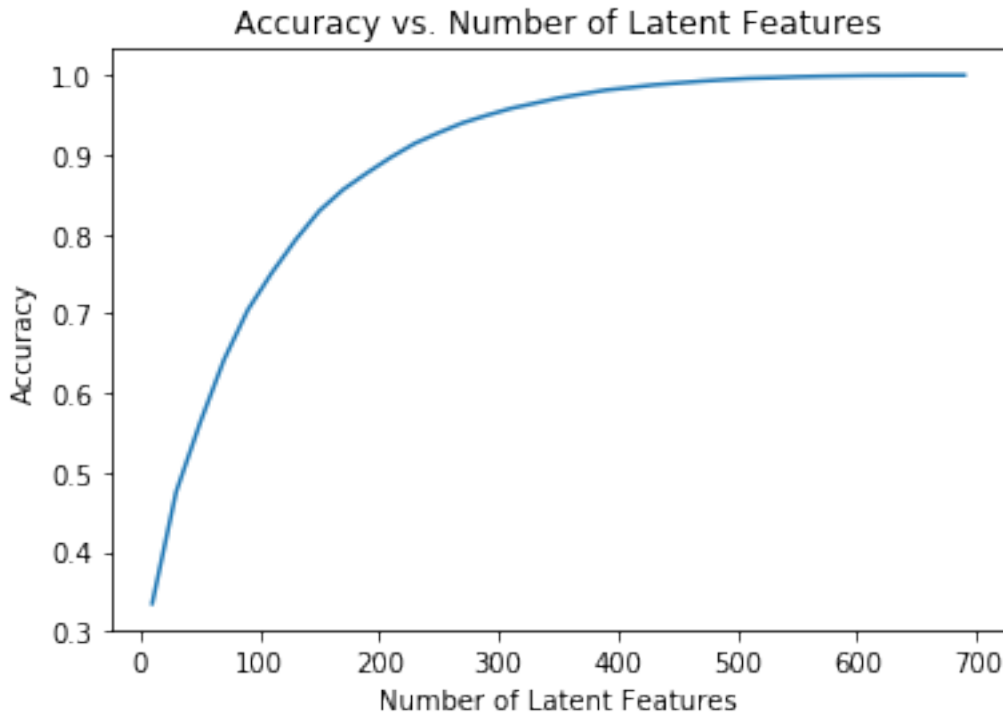
## Accuracy vs. Number of Latent Features



4. From the above, we can't really be sure how many features to use, because simply having a better way to predict the 1's and 0's of the matrix doesn't exactly give us an indication of if we are able to make good recommendations. Instead, we might split our dataset into a training and test set of data, as shown in the cell below.

Use the code from question 3 to understand the impact on accuracy of the training and test sets of data with different numbers of latent features. Using the split below:

- How many users can we make predictions for in the test set?

- How many users are we not able to make predictions for because of the cold start problem?
- How many articles can we make predictions for in the test set?

- How many articles are we not able to make predictions for because of the cold start problem?

```
In [120]: df_train = df.head(40000)
          df_test = df.tail(5993)

          def create_test_and_train_user_item(df_train, df_test):
              '''
              INPUT:
              df_train - training dataframe
              df_test - test dataframe

              OUTPUT:
              user_item_train - a user-item matrix of the training dataframe
```

```
                        (unique users for each row and unique articles for each column)
          user_item_test - a user-item matrix of the testing dataframe
                        (unique users for each row and unique articles for each column)
          test_idx - all of the test user ids
          test_arts - all of the test article ids


          '''
          # Your code here
          user_item_train=create_user_item_matrix(df_train)
          user_item_test=create_user_item_matrix(df_test)
          test_idx=user_item_test.index
          test_arts=user_item_test.columns

          return user_item_train, user_item_test, test_idx, test_arts

      user_item_train, user_item_test, test_idx, test_arts = create_test_and_train_user_item
```

In [122]: `user_item_train.head(5)`

Out[122]:

```
                 title                                                      \
article_id 0.0    2.0    4.0    8.0    9.0    12.0   14.0   15.0   16.0
user_id
1                0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0
2                0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0
3                0.0    0.0    0.0    0.0    0.0    1.0    0.0    0.0    0.0
4                0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0
5                0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0    0.0


                       ...                                                   \
article_id 18.0    ...   1434.0 1435.0 1436.0 1437.0 1439.0 1440.0 1441.0
user_id            ...
1                0.0   ...     0.0    0.0    1.0    0.0    1.0    0.0    0.0
2                0.0   ...     0.0    0.0    0.0    0.0    0.0    0.0    0.0
3                0.0   ...     0.0    0.0    1.0    0.0    0.0    0.0    0.0
4                0.0   ...     0.0    0.0    0.0    0.0    0.0    0.0    0.0
5                0.0   ...     0.0    0.0    0.0    0.0    0.0    0.0    0.0


article_id 1442.0 1443.0 1444.0
user_id
1                0.0    0.0    0.0
2                0.0    0.0    0.0
3                0.0    0.0    0.0
4                0.0    0.0    0.0
5                0.0    0.0    0.0

[5 rows x 714 columns]
```

In [124]: `print(u.shape, s.shape, vt.shape)`

23

```
(5149, 5149) (714,) (714, 714)
```

In [123]: *# Number of users in both sets*
          `len(user_item_test.index.intersection(user_item_train.index))`

Out[123]: 20

In [137]: *# movies in the test set are we not able to make predictions*
          `len(df_test.user_id.unique()) - len(np.intersect1d(df_train.user_id.unique(),df_test.u`

Out[137]: 662

In [138]: *# movies we can make predictions for in the test set*
          `len(np.intersect1d(df_train.article_id.unique(),df_test.article_id.unique()))`

Out[138]: 574

In [139]: *# users in the test set are we not able to make predictions*
          `len(df_test.article_id.unique()) - len(np.intersect1d(df_train.article_id.unique(),df_`

Out[139]: 0

In [136]: *# Replace the values in the dictionary below*
          `a = 662`
          `b = 574`
          `c = 20`
          `d = 0`


          `sol_4_dict = {`
              `'How many users can we make predictions for in the test set?':c, # letter here,`
              `'How many users in the test set are we not able to make predictions for because of`
              `'How many movies can we make predictions for in the test set?': b, # letter here,`
              `'How many movies in the test set are we not able to make predictions for because o`
          `}`

          `t.sol_4_test(sol_4_dict)`

```
Awesome job!  That's right!  All of the test movies are in the training data, but there are only
```

5. Now use the **user_item_train** dataset from above to find U, S, and V transpose using SVD. Then find the subset of rows in the **user_item_test** dataset that you can predict using this matrix decomposition with different numbers of latent features to see how many features makes sense to keep based on the accuracy on the test data. This will require combining what was done in questions 2 - 4.

Use the cells below to explore how well SVD works towards making predictions for recommendations on the test data.

```
In [140]: # fit SVD on the user_item_train matrix
          u_train, s_train, vt_train = np.linalg.svd(user_item_train) # fit svd similar to above

In [141]: u_train.shape, s_train.shape, vt_train.shape

Out[141]: ((4487, 4487), (714,), (714, 714))

In [ ]: # Use these cells to see how well you can use the training
        # decomposition to predict on test data

In [149]: num_latent_feats = np.arange(10,700+10,20)
          sum_errs_train = []
          sum_errs_test = []
          user_item_test = user_item_test.loc[user_item_test.index.isin(user_item_train.index),
          u_test = u_train[user_item_train.index.isin(user_item_test.index), :]
          vt_test = vt_train[:, user_item_train.columns.isin(test_arts)]


          for k in num_latent_feats:
              # restructure with k latent features
              s_new_train, u_new_train, vt_new_train = np.diag(s_train[:k]), u_train[:, :k], vt_


              s_new_test, u_new_test, vt_new_test = s_new_train, u_test[:, :k], vt_test[:k, :]

              # take dot product
              user_item_est_train = np.around(np.dot(np.dot(u_new_train, s_new_train), vt_new_tr
              user_item_est_test = np.around(np.dot(np.dot(u_new_test, s_new_test), vt_new_test)

              # compute error for each prediction to actual value
              diffs_train = np.subtract(user_item_train, user_item_est_train)
              diffs_test = np.subtract(user_item_test, user_item_est_test)

              # total errors and keep track of them
              err_train = np.sum(np.sum(np.abs(diffs_train)))
              err_test = np.sum(np.sum(np.abs(diffs_test)))

              sum_errs_train.append(err_train)

              sum_errs_test.append(err_test)

In [150]: fig, ax1 = plt.subplots()
          ax2 = ax1.twinx()

          ax1.plot(num_latent_feats, 1 - np.array(sum_errs_train)/df.shape[0], label="Train accu
          ax2.plot(num_latent_feats, 1 - np.array(sum_errs_test)/df.shape[0], color='green', lab

          handler1, label1 = ax1.get_legend_handles_labels()
          handler2, label2 = ax2.get_legend_handles_labels()
```
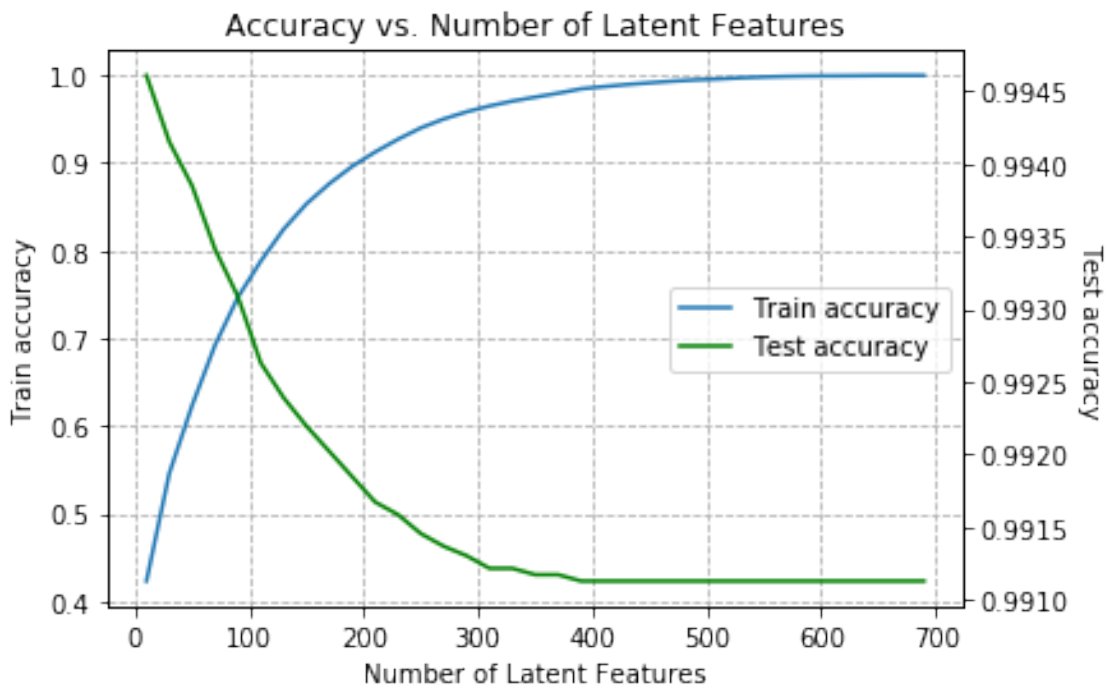
```
ax1.legend(handler1+handler2, label1+label2, loc='center right')

ax1.set_title('Accuracy vs. Number of Latent Features')
ax1.grid(linestyle='--')

ax1.set_xlabel('Number of Latent Features')
ax1.set_ylabel('Train accuracy')
ax2.set_ylabel('Test accuracy', rotation=270, labelpad=12)

plt.show()
```



6. Use the cell below to comment on the results you found in the previous question. Given the circumstances of your results, discuss what you might do to determine if the recommendations you make with any of the above recommendation systems are an improvement to how users currently find articles?

**Your response here.**

From the plot above we can see that as the number of latent features increases, the train acccuracy increases but the test accuracy decrease. This is a clear case of overfitting. Our model overfits the training dataset indicating that the more latent features, the more overfitting will happen. Based on this I would try to keep fewer latent features.

However, we only have data for 20 overlapping users, I believe this number is too small for any statistical significance as it creates difficulties to predict the accuracy with the limited number of users. Based on this, I cannot detrmine with high certainty that the SVD recommendations work well in this case.

Furthermore, we could use other recommendation methods to improve our recommendation, like collaborative filtering or content based recommendation. Then we could use A/B testing to check which model actually works well in practice.

### Extras Using your workbook, you could now save your recommendations for each user, develop a class to make new predictions and update your results, and make a flask app to deploy your results. These tasks are beyond what is required for this project. However, from what you learned in the lessons, you certainly capable of taking these tasks on to improve upon your work here!

## 1.2 Conclusion

Congratulations! You have reached the end of the Recommendations with IBM project!

## 1.3 Directions to Submit

Before you submit your project, you need to create a .html or .pdf version of this notebook in the workspace here. To do that, run the code cell below. If it worked correctly, you should get a return code of 0, and you should see the generated .html file in the workspace directory (click on the orange Jupyter icon in the upper left).

Alternatively, you can download this report as .html via the **File** > **Download as** submenu, and then manually upload it into the workspace directory by clicking on the orange Jupyter icon in the upper left, then using the Upload button.

Once you've done this, you can submit your project by clicking on the "Submit Project" button in the lower right here. This will create and submit a zip file with this .ipynb doc and the .html or .pdf version you created. Congratulations!

```
In [152]: from subprocess import call
          call(['python', '-m', 'nbconvert', 'Recommendations_with_IBM.ipynb'])

Out[152]: 0

In [ ]:
```