In [2]:
```python
import pandas as pd
import numpy as np
from scipy.sparse import csr_matrix
from sklearn.neighbors import NearestNeighbors
# import matplotlib.pyplot as plt
# import seaborn as sns
movies = pd.read_csv(r"C:\Users\chinmayee\Downloads\movies.csv")
ratings = pd.read_csv(r"C:\Users\chinmayee\Downloads\ratings (1).csv")

print(movies.head())
```

```
   movieId                               title  \
0        1                    Toy Story (1995)
1        2                      Jumanji (1995)
2        3             Grumpier Old Men (1995)
3        4            Waiting to Exhale (1995)
4        5  Father of the Bride Part II (1995)

                                        genres
0  Adventure|Animation|Children|Comedy|Fantasy
1                   Adventure|Children|Fantasy
2                               Comedy|Romance
3                         Comedy|Drama|Romance
4                                       Comedy
```

In [4]:
```python
ratings.head()
```

Out[4]:

|   | userId | movieId | rating | timestamp |
|---|--------|---------|--------|-----------|
| 0 | 1 | 1 | 4.0 | 964982703 |
| 1 | 1 | 3 | 4.0 | 964981247 |
| 2 | 1 | 6 | 4.0 | 964982224 |
| 3 | 1 | 47 | 5.0 | 964983815 |
| 4 | 1 | 50 | 5.0 | 964982931 |

In [5]:
```python
final_dataset = ratings.pivot(index='movieId',columns='userId',values='rating')
final_dataset.head()
```

Out[5]:

| userId | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... | 601 | 602 | 603 | 604 |
|--------|---|---|---|---|---|---|---|---|---|-----|-----|-----|-----|-----|-----|
| **movieId** | | | | | | | | | | | | | | | |
| 1 | 4.0 | NaN | NaN | NaN | 4.0 | NaN | 4.5 | NaN | NaN | NaN | ... | 4.0 | NaN | 4.0 | 3.0 |
| 2 | NaN | NaN | NaN | NaN | NaN | 4.0 | NaN | 4.0 | NaN | NaN | ... | NaN | 4.0 | NaN | 5.0 |
| 3 | 4.0 | NaN | NaN | NaN | NaN | 5.0 | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 4 | NaN | NaN | NaN | NaN | NaN | 3.0 | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN |
| 5 | NaN | NaN | NaN | NaN | NaN | 5.0 | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | 3.0 |

5 rows × 610 columns

In [6]:
```python
final_dataset.fillna(0,inplace=True)
final_dataset.head()
```

Out[6]:

| userId | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... | 601 | 602 | 603 | 604 | 605 | 606 | 607 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| movieId | | | | | | | | | | | | | | | | | | |
| 1 | 4.0 | 0.0 | 0.0 | 0.0 | 4.0 | 0.0 | 4.5 | 0.0 | 0.0 | 0.0 | ... | 4.0 | 0.0 | 4.0 | 3.0 | 4.0 | 2.5 | 4.0 |
| 2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 4.0 | 0.0 | 4.0 | 0.0 | 0.0 | ... | 0.0 | 4.0 | 0.0 | 5.0 | 3.5 | 0.0 | 0.0 |
| 3 | 4.0 | 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 | 0.0 | 0.0 | 0.0 |
| 4 | 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 | 0.0 | 0.0 | 0.0 |
| 5 | 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 | 3.0 | 0.0 | 0.0 | 0.0 |

5 rows × 610 columns

In [7]:
```python
no_user_voted = ratings.groupby('movieId')['rating'].agg('count')
no_movies_voted = ratings.groupby('userId')['rating'].agg('count')
```
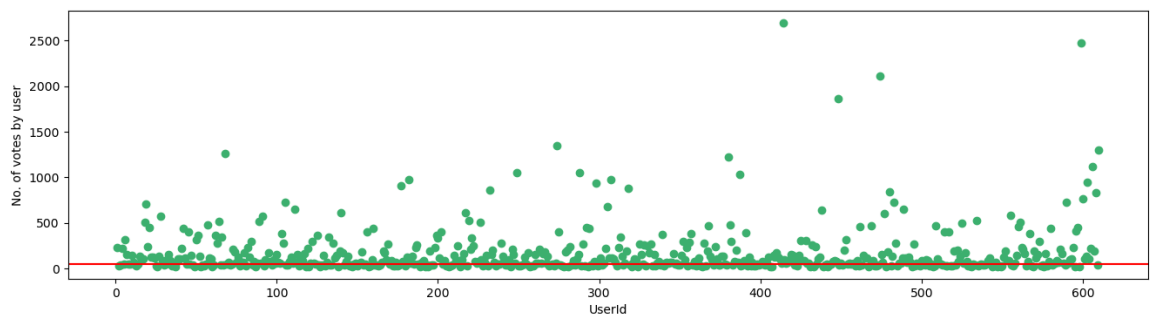
In [10]:
```python
#Let's visualize the number of users who voted with our threshold of 10.
import matplotlib.pyplot as plt
f,ax = plt.subplots(1,1,figsize=(16,4))
# ratings['rating'].plot(kind='hist')
plt.scatter(no_user_voted.index,no_user_voted,color='mediumseagreen')
plt.axhline(y=10,color='r')
plt.xlabel('MovieId')
plt.ylabel('No. of users voted')
plt.show()
```



In [11]:
```python
final_dataset = final_dataset.loc[no_user_voted[no_user_voted > 10].index,:]
```

In [12]:
```python
#Let's visualize the number of votes by each user with our threshold of 50.

f,ax = plt.subplots(1,1,figsize=(16,4))
plt.scatter(no_movies_voted.index,no_movies_voted,color='mediumseagreen')
plt.axhline(y=50,color='r')
plt.xlabel('UserId')
plt.ylabel('No. of votes by user')
plt.show()
```

```
In [13]:  final_dataset=final_dataset.loc[:,no_movies_voted[no_movies_voted > 50].index]
          final_dataset
```

Out[13]:

| userId movieId | 1 | 4 | 6 | 7 | 10 | 11 | 15 | 16 | 17 | 18 | ... | 600 | 601 | 602 | 603 | 604 | 605 | 606 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 4.0 | 0.0 | 0.0 | 4.5 | 0.0 | 0.0 | 2.5 | 0.0 | 4.5 | 3.5 | ... | 2.5 | 4.0 | 0.0 | 4.0 | 3.0 | 4.0 | 2.5 |
| **2** | 0.0 | 0.0 | 4.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 3.0 | ... | 4.0 | 0.0 | 4.0 | 0.0 | 5.0 | 3.5 | 0.0 |
| **3** | 4.0 | 0.0 | 5.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **5** | 0.0 | 0.0 | 5.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 2.5 | 0.0 | 0.0 | 0.0 | 3.0 | 0.0 | 0.0 |
| **6** | 4.0 | 0.0 | 4.0 | 0.0 | 0.0 | 5.0 | 0.0 | 0.0 | 0.0 | 4.0 | ... | 0.0 | 0.0 | 3.0 | 4.0 | 3.0 | 0.0 | 0.0 |
| **...** | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | .. |
| **174055** | 0.0 | 0.0 | 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 |
| **176371** | 0.0 | 0.0 | 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 |
| **177765** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 4.5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **179819** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **187593** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

2121 rows × 378 columns

```
In [15]:  sample = np.array([[0,0,3,0,0],[4,0,0,0,2],[0,0,0,0,1]])
          sparsity = 1.0 - ( np.count_nonzero(sample) / float(sample.size) )
          print(sparsity)
```

```
0.7333333333333334
```

```
In [16]:  csr_sample = csr_matrix(sample)
          print(csr_sample)
```

```
  (0, 2)        3
  (1, 0)        4
  (1, 4)        2
  (2, 4)        1
```

```
In [17]:  csr_data = csr_matrix(final_dataset.values)
          final_dataset.reset_index(inplace=True)
```

```
In [18]:  knn = NearestNeighbors(metric='cosine', algorithm='brute', n_neighbors=20, n_job
          knn.fit(csr_data)
```

Out[18]:

```
▼                          NearestNeighbors

NearestNeighbors(algorithm='brute', metric='cosine', n_jobs=-1, n_neigh
bors=20)
```

In [19]:

```python
def get_movie_recommendation(movie_name):
    n_movies_to_reccomend = 10
    movie_list = movies[movies['title'].str.contains(movie_name)]
    if len(movie_list):
        movie_idx= movie_list.iloc[0]['movieId']
        movie_idx = final_dataset[final_dataset['movieId'] == movie_idx].index[0
        distances , indices = knn.kneighbors(csr_data[movie_idx],n_neighbors=n_m
        rec_movie_indices = sorted(list(zip(indices.squeeze().tolist(),distances
        recommend_frame = []
        for val in rec_movie_indices:
            movie_idx = final_dataset.iloc[val[0]]['movieId']
            idx = movies[movies['movieId'] == movie_idx].index
            recommend_frame.append({'Title':movies.iloc[idx]['title'].values[0],
        df = pd.DataFrame(recommend_frame,index=range(1,n_movies_to_reccomend+1)
        return df
    else:
        return "No movies found. Please check your input"
```

In [20]:

```python
get_movie_recommendation('Iron Man')
```

Out[20]:

| | Title | Distance |
|---|---|---|
| 1 | Up (2009) | 0.368857 |
| 2 | Guardians of the Galaxy (2014) | 0.368758 |
| 3 | Watchmen (2009) | 0.368558 |
| 4 | Star Trek (2009) | 0.366029 |
| 5 | Batman Begins (2005) | 0.362759 |
| 6 | Avatar (2009) | 0.310893 |
| 7 | Iron Man 2 (2010) | 0.307492 |
| 8 | WALL·E (2008) | 0.298138 |
| 9 | Dark Knight, The (2008) | 0.285835 |
| 10 | Avengers, The (2012) | 0.285319 |

In [21]:

```python
get_movie_recommendation('Memento')
```

Out[21]:

| | Title | Distance |
|---|---|---|
| **1** | American Beauty (1999) | 0.389346 |
| **2** | American History X (1998) | 0.388615 |
| **3** | Pulp Fiction (1994) | 0.386235 |
| **4** | Lord of the Rings: The Return of the King, The... | 0.371622 |
| **5** | Kill Bill: Vol. 1 (2003) | 0.350167 |
| **6** | Lord of the Rings: The Two Towers, The (2002) | 0.348358 |
| **7** | Eternal Sunshine of the Spotless Mind (2004) | 0.346196 |
| **8** | Matrix, The (1999) | 0.326215 |
| **9** | Lord of the Rings: The Fellowship of the Ring,... | 0.316777 |
| **10** | Fight Club (1999) | 0.272380 |

In [ ]: