

```
In [4]: import pandas as pd
```

```
In [6]: movies = pd.read_csv(r"C:\Users\smak_\Desktop\NareshITechnologies\EDA-Aug25th\my-wo
```

```
In [7]: movies
```

```
Out[7]:
```

	Film	Genre	Rotten Tomatoes Ratings %	Audience Ratings %	Budget (million \$)	Year of release
0	(500) Days of Summer	Comedy	87	81	8	2009
1	10,000 B.C.	Adventure	9	44	105	2008
2	12 Rounds	Action	30	52	20	2009
3	127 Hours	Adventure	93	84	18	2010
4	17 Again	Comedy	55	70	20	2009
...
554	Your Highness	Comedy	26	36	50	2011
555	Youth in Revolt	Comedy	68	52	18	2009
556	Zodiac	Thriller	89	73	65	2007
557	Zombieland	Action	90	87	24	2009
558	Zookeeper	Comedy	14	42	80	2011

559 rows × 6 columns

```
In [9]: type(movies)
```

```
Out[9]: pandas.core.frame.DataFrame
```

```
In [10]: movies.isna()
```

Out[10]:

	Film	Genre	Rotten Tomatoes Ratings %	Audience Ratings %	Budget (million \$)	Year of release
0	False	False	False	False	False	False
1	False	False	False	False	False	False
2	False	False	False	False	False	False
3	False	False	False	False	False	False
4	False	False	False	False	False	False
...
554	False	False	False	False	False	False
555	False	False	False	False	False	False
556	False	False	False	False	False	False
557	False	False	False	False	False	False
558	False	False	False	False	False	False

559 rows × 6 columns

In [11]:

movies.isnull().sum()

Out[11]:

Film0
Genre0
Rotten Tomatoes Ratings %0
Audience Ratings %0
Budget (million \$)0
Year of release0
dtype: int64

In [14]:

len(movies)

Out[14]:

559

In [17]:

movies.columns

Out[17]:

Index(['Film', 'Genre', 'Rotten Tomatoes Ratings %', 'Audience Ratings %',
 'Budget (million \$)', 'Year of release'],
 dtype='object')

In [18]:

movies.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 559 entries, 0 to 558
Data columns (total 6 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Film                                559 non-null    object
1   Genre                              559 non-null    object
2   Rotten Tomatoes Ratings %          559 non-null    int64
3   Audience Ratings %                 559 non-null    int64
4   Budget (million $)                 559 non-null    int64
5   Year of release                     559 non-null    int64
dtypes: int64(4), object(2)
memory usage: 26.3+ KB
```

In [20]: `movies.shape`

Out[20]: (559, 6)

In [21]: `movies.head()`

Out[21]:

	Film	Genre	Rotten Tomatoes Ratings %	Audience Ratings %	Budget (million \$)	Year of release
0	(500) Days of Summer	Comedy	87	81	8	2009
1	10,000 B.C.	Adventure	9	44	105	2008
2	12 Rounds	Action	30	52	20	2009
3	127 Hours	Adventure	93	84	18	2010
4	17 Again	Comedy	55	70	20	2009

In [22]: `movies.tail()`

Out[22]:

	Film	Genre	Rotten Tomatoes Ratings %	Audience Ratings %	Budget (million \$)	Year of release
554	Your Highness	Comedy	26	36	50	2011
555	Youth in Revolt	Comedy	68	52	18	2009
556	Zodiac	Thriller	89	73	65	2007
557	Zombieland	Action	90	87	24	2009
558	Zookeeper	Comedy	14	42	80	2011

In [23]: `movies.columns`

```
Out[23]: Index(['Film', 'Genre', 'Rotten Tomatoes Ratings %', 'Audience Ratings %',
              'Budget (million $)', 'Year of release'],
              dtype='object')
```

```
In [24]: movies.columns = ['Film', 'Genre', 'CriticRating', 'AudienceRating', 'BudgetMillions']
```

```
In [27]: movies.head(1) # removed spaces and removed noise characters
```

```
Out[27]:
```

	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
0	(500) Days of Summer	Comedy	87	81	8	2009

```
In [30]: movies.describe() # Descriptive statistics
```

```
Out[30]:
```

	CriticRating	AudienceRating	BudgetMillions	Year
count	559.000000	559.000000	559.000000	559.000000
mean	47.309481	58.744186	50.236136	2009.152057
std	26.413091	16.826887	48.731817	1.362632
min	0.000000	0.000000	0.000000	2007.000000
25%	25.000000	47.000000	20.000000	2008.000000
50%	46.000000	58.000000	35.000000	2009.000000
75%	70.000000	72.000000	65.000000	2010.000000
max	97.000000	96.000000	300.000000	2011.000000

```
In [31]: movies.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 559 entries, 0 to 558
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Film            559 non-null    object
1   Genre           559 non-null    object
2   CriticRating    559 non-null    int64
3   AudienceRating  559 non-null    int64
4   BudgetMillions  559 non-null    int64
5   Year            559 non-null    int64
dtypes: int64(4), object(2)
memory usage: 26.3+ KB
```

```
In [32]: movies.Film = movies.Film.astype('category')
```

```
In [33]: movies.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 559 entries, 0 to 558
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Film                  559 non-null    category
1   Genre                  559 non-null    object
2   CriticRating           559 non-null    int64
3   AudienceRating         559 non-null    int64
4   BudgetMillions         559 non-null    int64
5   Year                   559 non-null    int64
dtypes: category(1), int64(4), object(1)
memory usage: 43.6+ KB
```

In [34]: `movies.describe()`

Out[34]:

	CriticRating	AudienceRating	BudgetMillions	Year
count	559.000000	559.000000	559.000000	559.000000
mean	47.309481	58.744186	50.236136	2009.152057
std	26.413091	16.826887	48.731817	1.362632
min	0.000000	0.000000	0.000000	2007.000000
25%	25.000000	47.000000	20.000000	2008.000000
50%	46.000000	58.000000	35.000000	2009.000000
75%	70.000000	72.000000	65.000000	2010.000000
max	97.000000	96.000000	300.000000	2011.000000

In [36]: `movies.Genre = movies.Genre.astype('category')`
`movies.Year = movies.Year.astype('category')`

In [39]: `movies.Genre`

Out[39]:

0	Comedy
1	Adventure
2	Action
3	Adventure
4	Comedy
...	
554	Comedy
555	Comedy
556	Thriller
557	Action
558	Comedy

Name: Genre, Length: 559, dtype: category
Categories (7, object): ['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance', 'Thriller']

In [40]: `movies.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 559 entries, 0 to 558
Data columns (total 6 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Film                  559 non-null    category
1   Genre                  559 non-null    category
2   CriticRating           559 non-null    int64
3   AudienceRating         559 non-null    int64
4   BudgetMillions         559 non-null    int64
5   Year                   559 non-null    category
dtypes: category(3), int64(3)
memory usage: 36.5 KB
```

In [42]: `movies.describe()`

Out[42]:

	CriticRating	AudienceRating	BudgetMillions
count	559.000000	559.000000	559.000000
mean	47.309481	58.744186	50.236136
std	26.413091	16.826887	48.731817
min	0.000000	0.000000	0.000000
25%	25.000000	47.000000	20.000000
50%	46.000000	58.000000	35.000000
75%	70.000000	72.000000	65.000000
max	97.000000	96.000000	300.000000

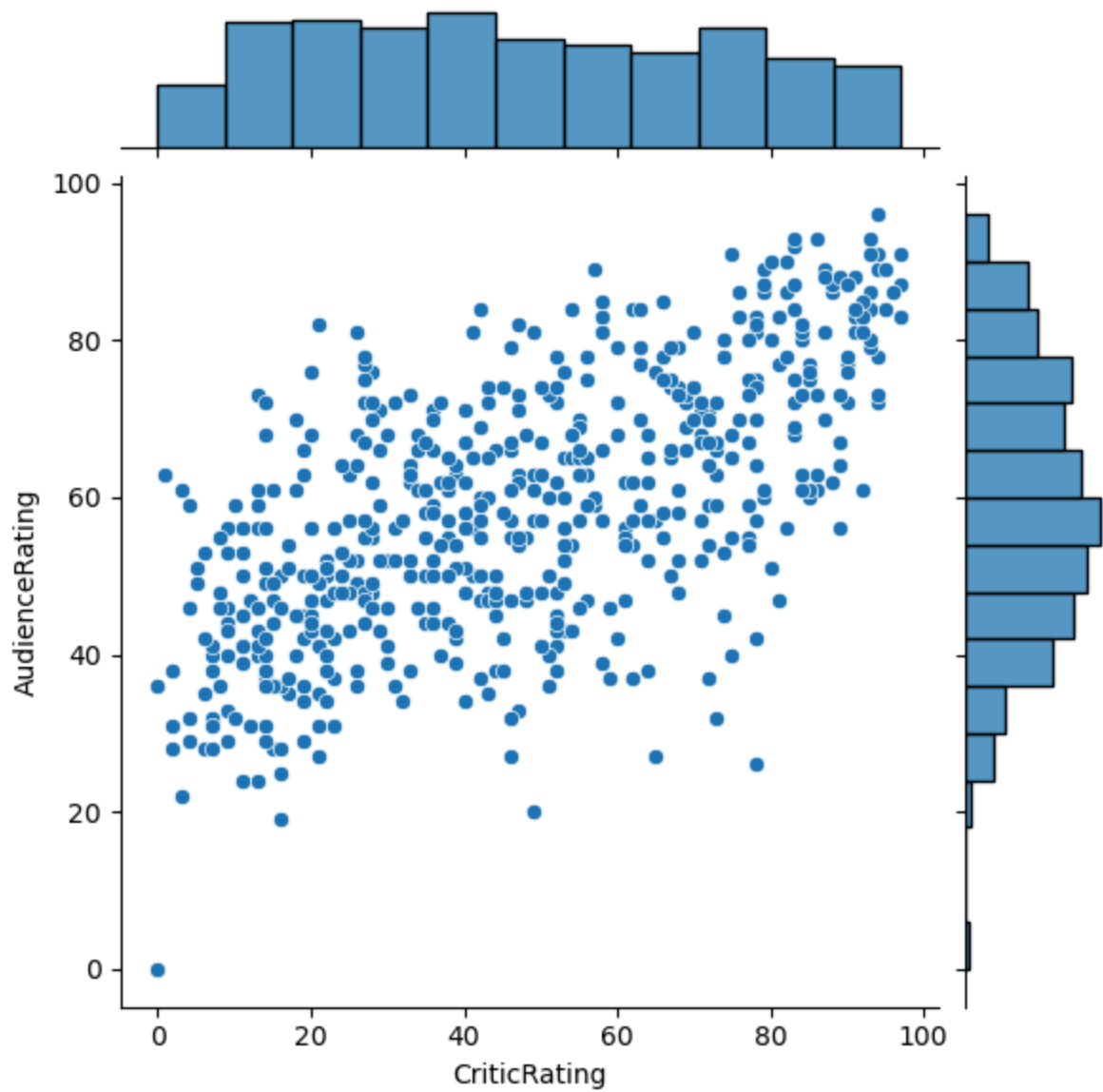
In [44]: *# How to work with joint plots*

```
from matplotlib import pyplot as plt # For visualization
import seaborn as sns # Advanced visualization

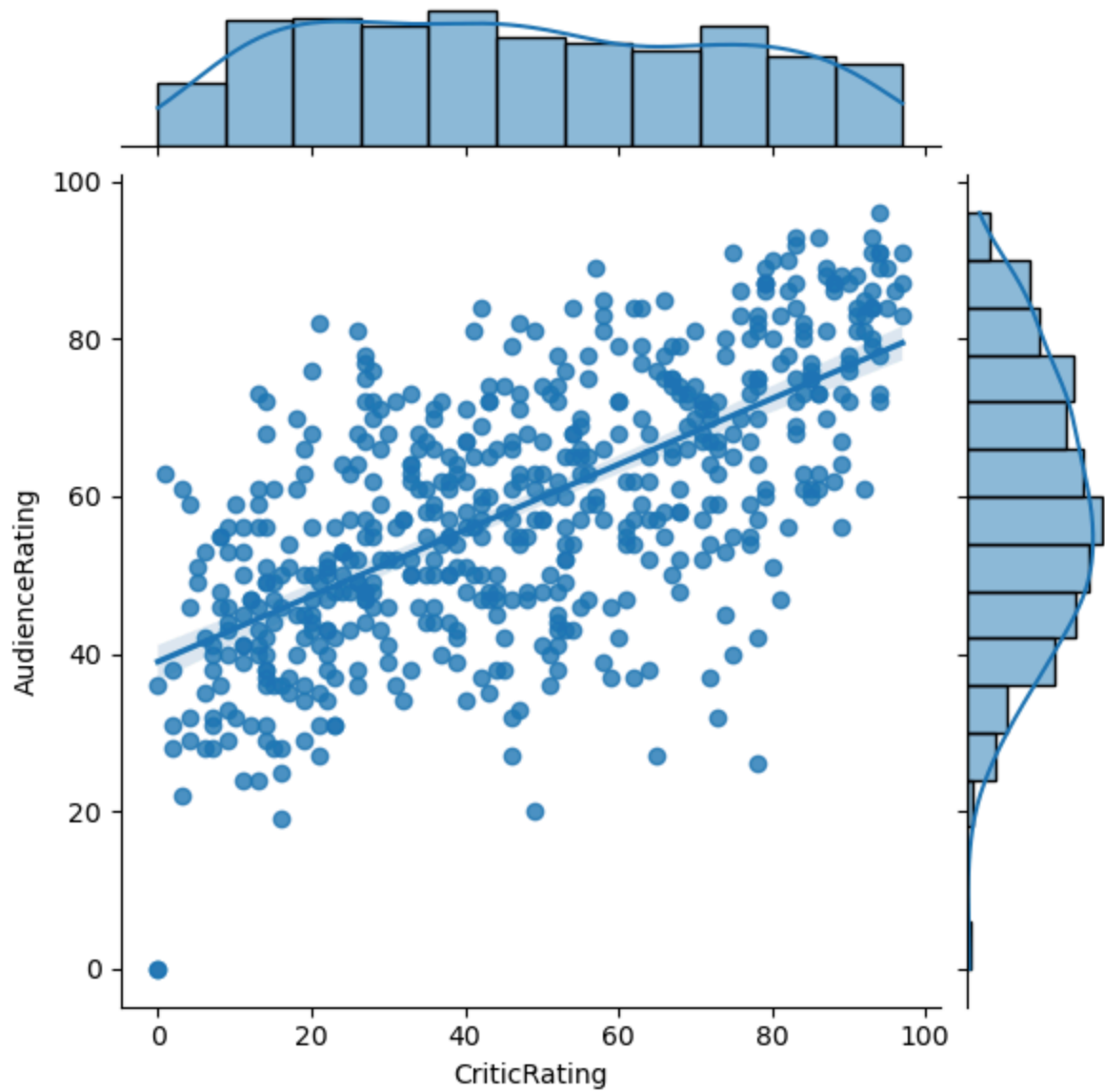
%%matplotlib inline # All the plot should be inside the line

import warnings
warnings.filterwarnings('ignore')
```

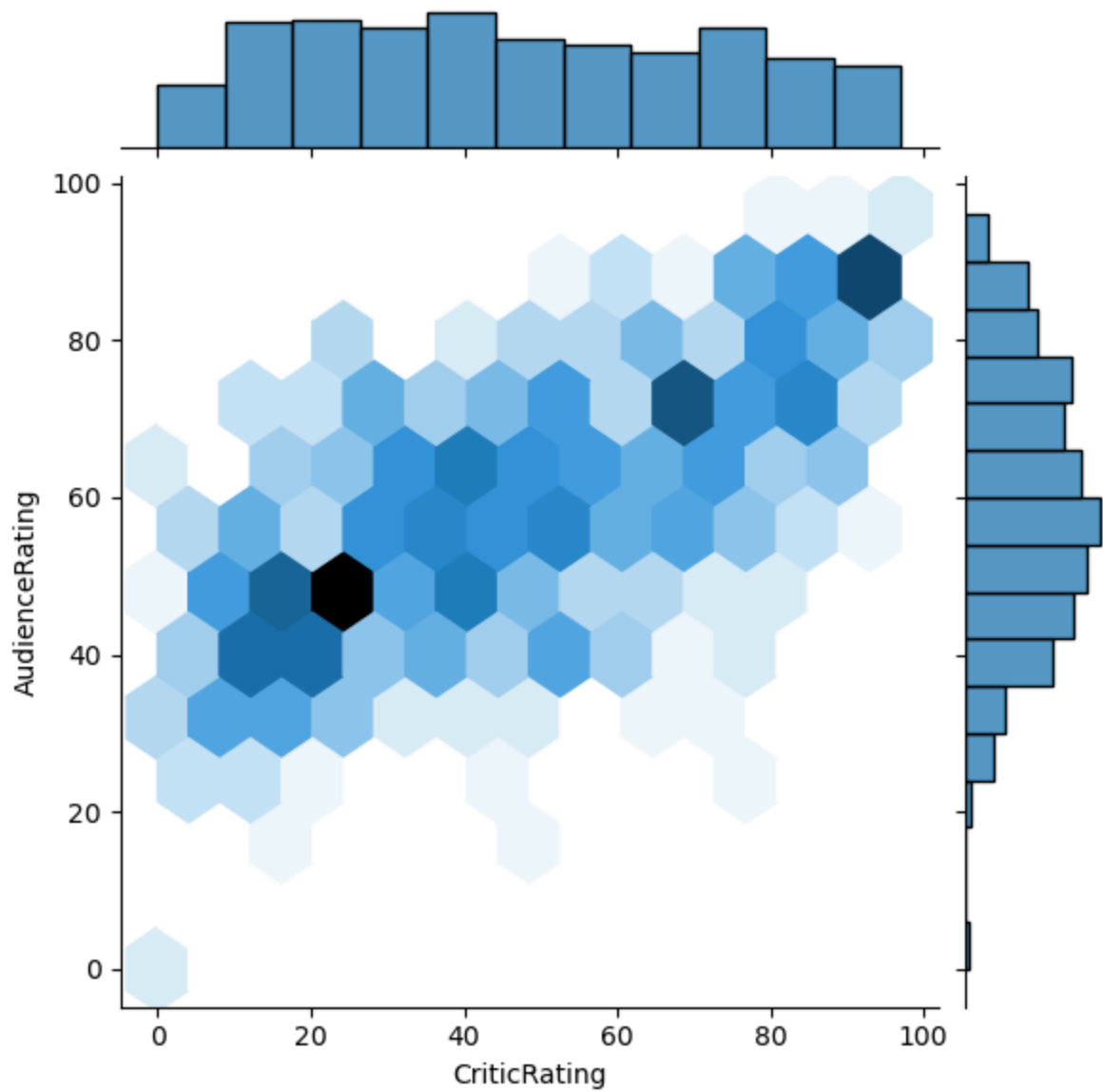
In [47]: `j = sns.jointplot(data = movies, x = 'CriticRating', y = 'AudienceRating', kind = 's`



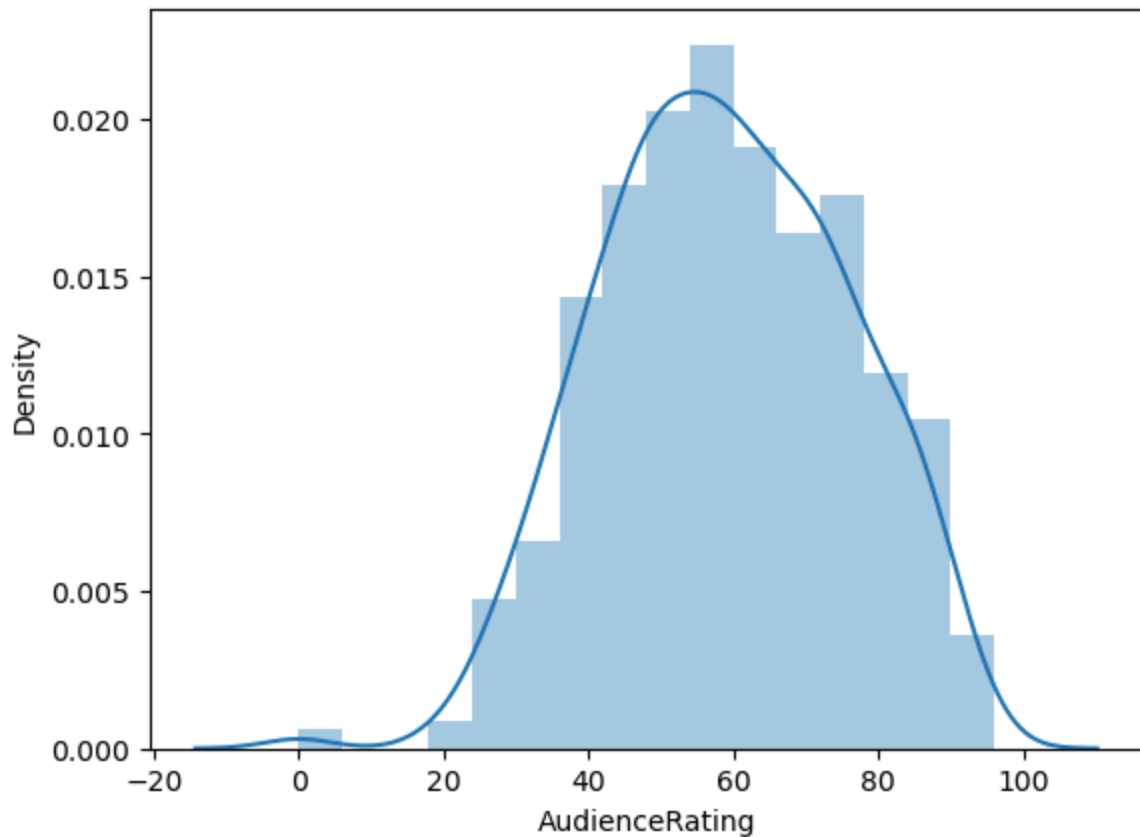
```
In [48]: j = sns.jointplot(data = movies, x = 'CriticRating', y = 'AudienceRating', kind = 'r
```



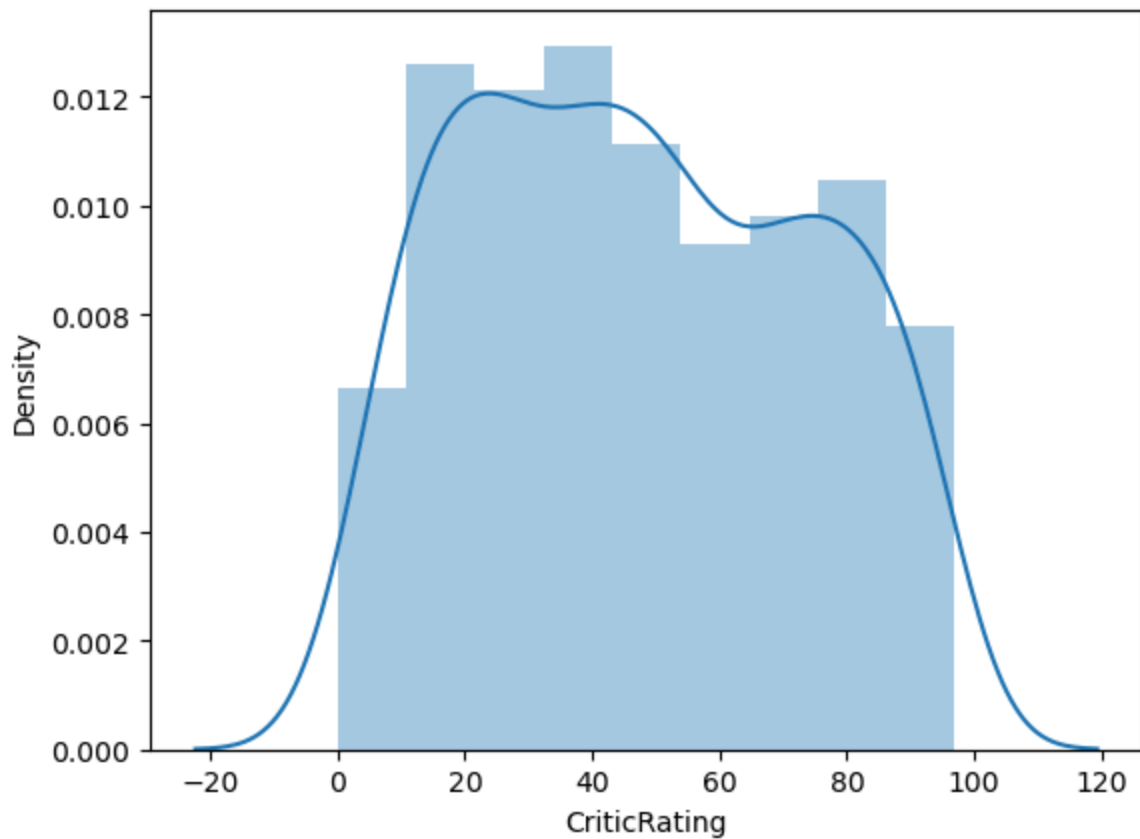
```
In [49]: j = sns.jointplot(data = movies, x = 'CriticRating', y = 'AudienceRating', kind = 'h
```

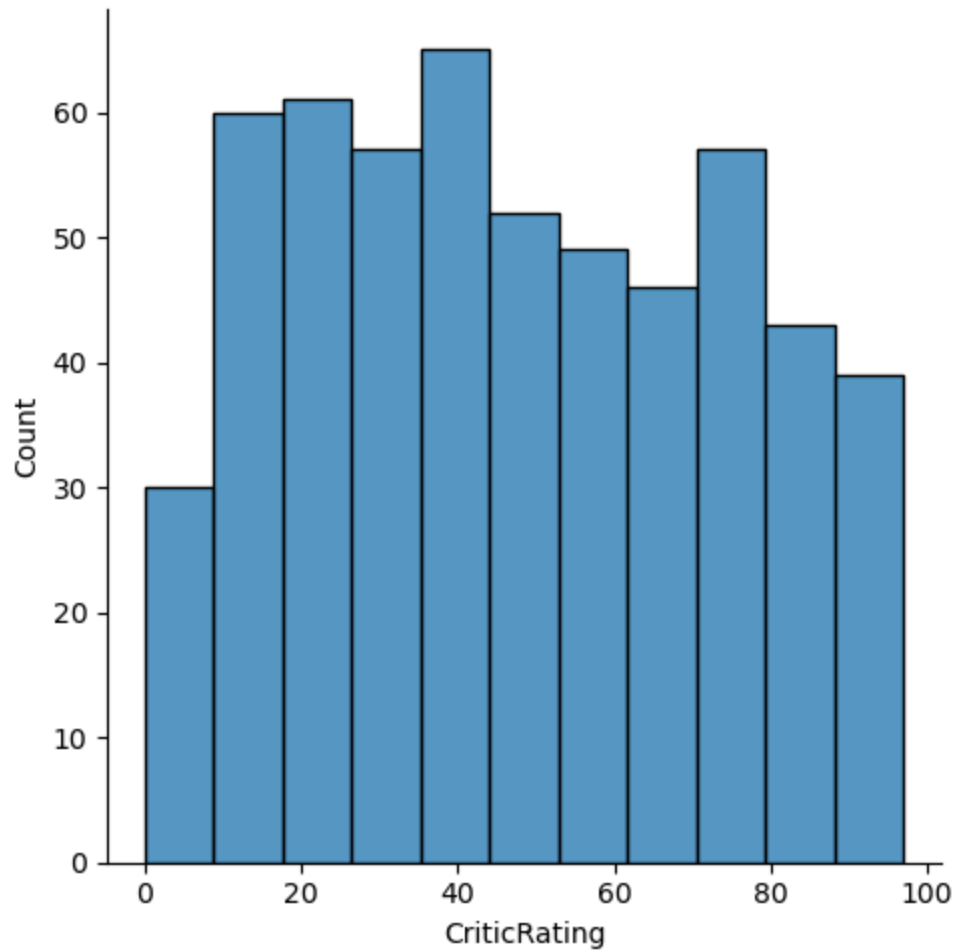
```
In [50]: m1 = sns.distplot(movies.AudienceRating)
```



```
In [51]: m1 = sns.distplot(movies.CriticRating)
```

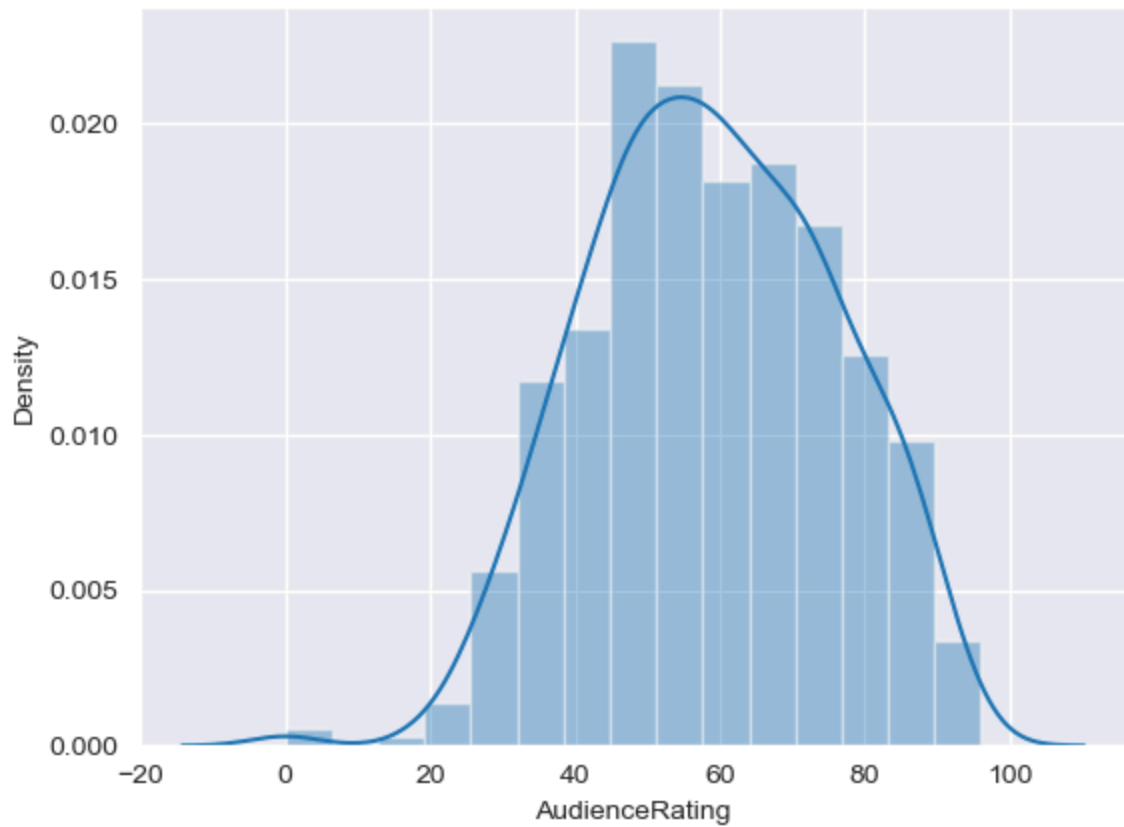


```
In [52]: m1 = sns.displot(movies.CriticRating)
```



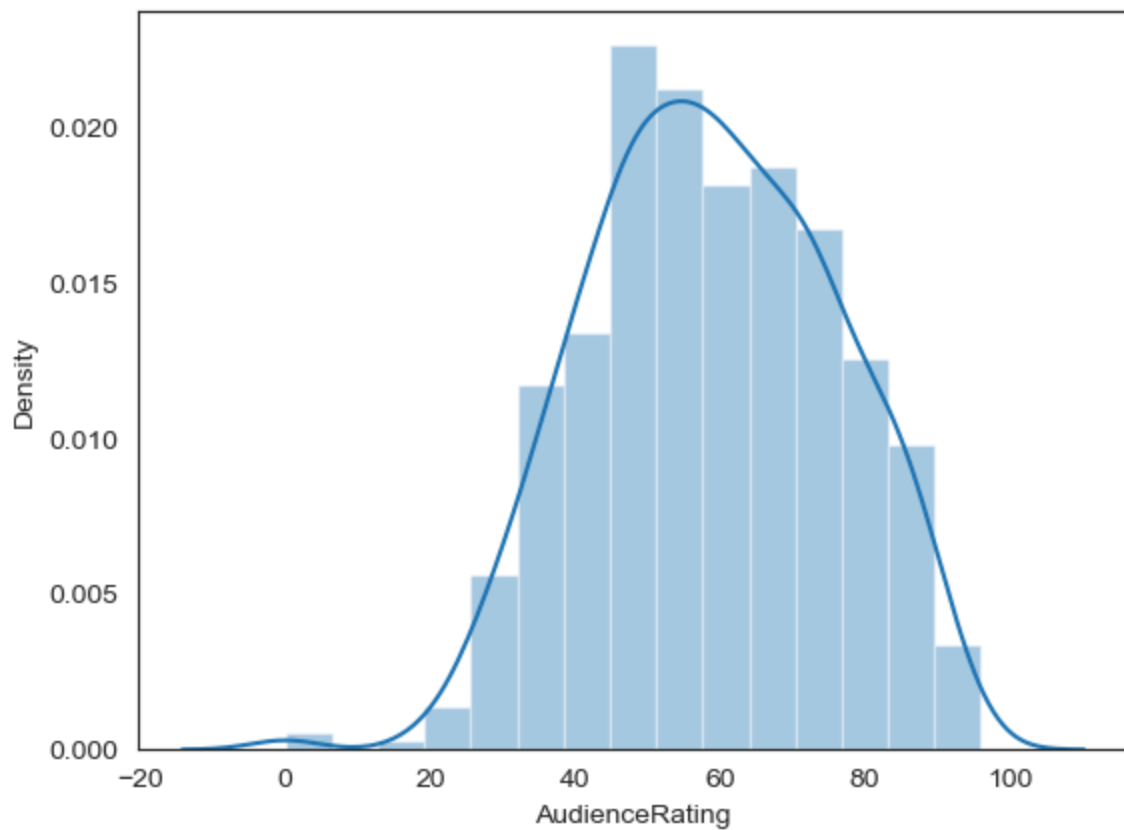
```
In [53]: sns.set_style('darkgrid')
```

```
In [54]: m2 = sns.distplot(movies.AudienceRating, bins = 15)
```

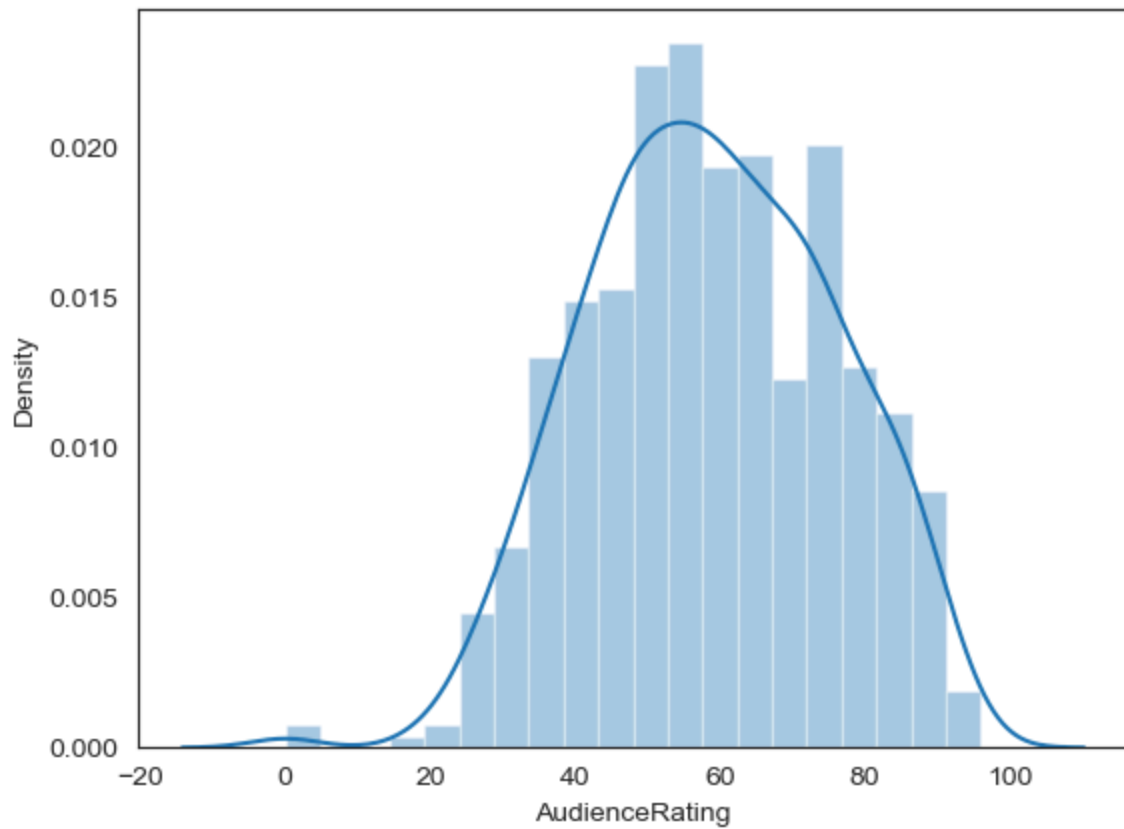


```
In [55]: sns.set_style('white')
```

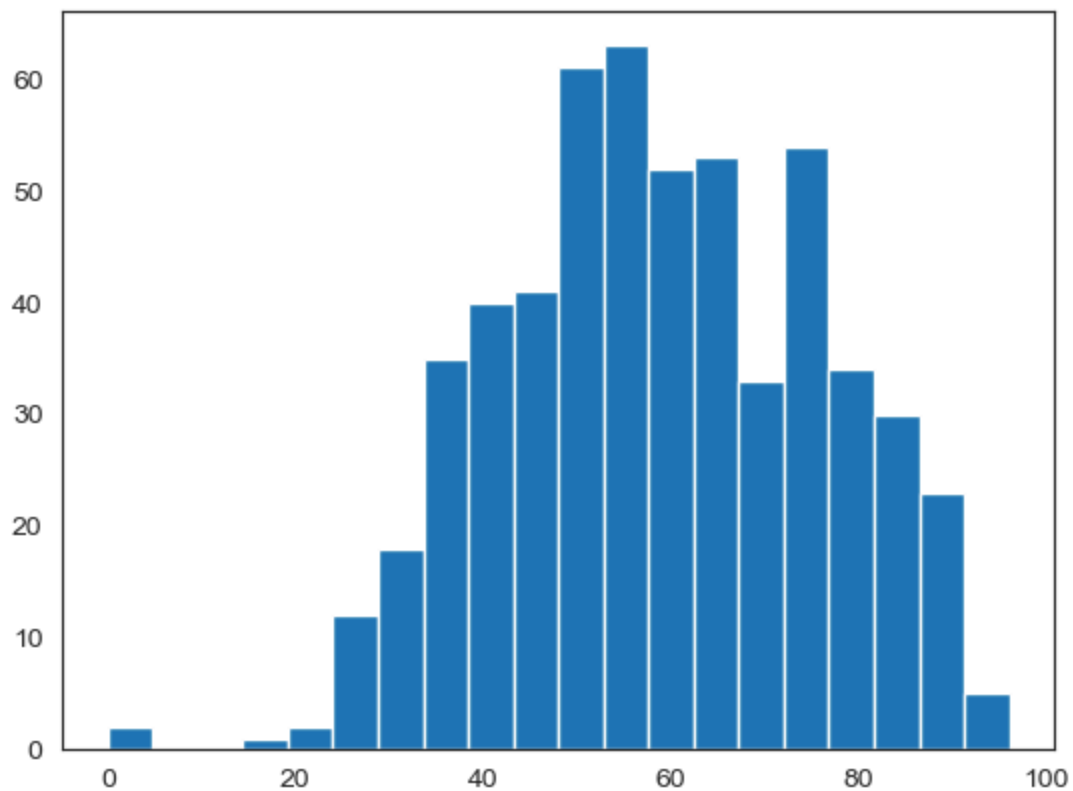
```
In [56]: m2 = sns.distplot(movies.AudienceRating, bins = 15)
```



```
In [57]: m2 = sns.distplot(movies.AudienceRating, bins = 20)
```

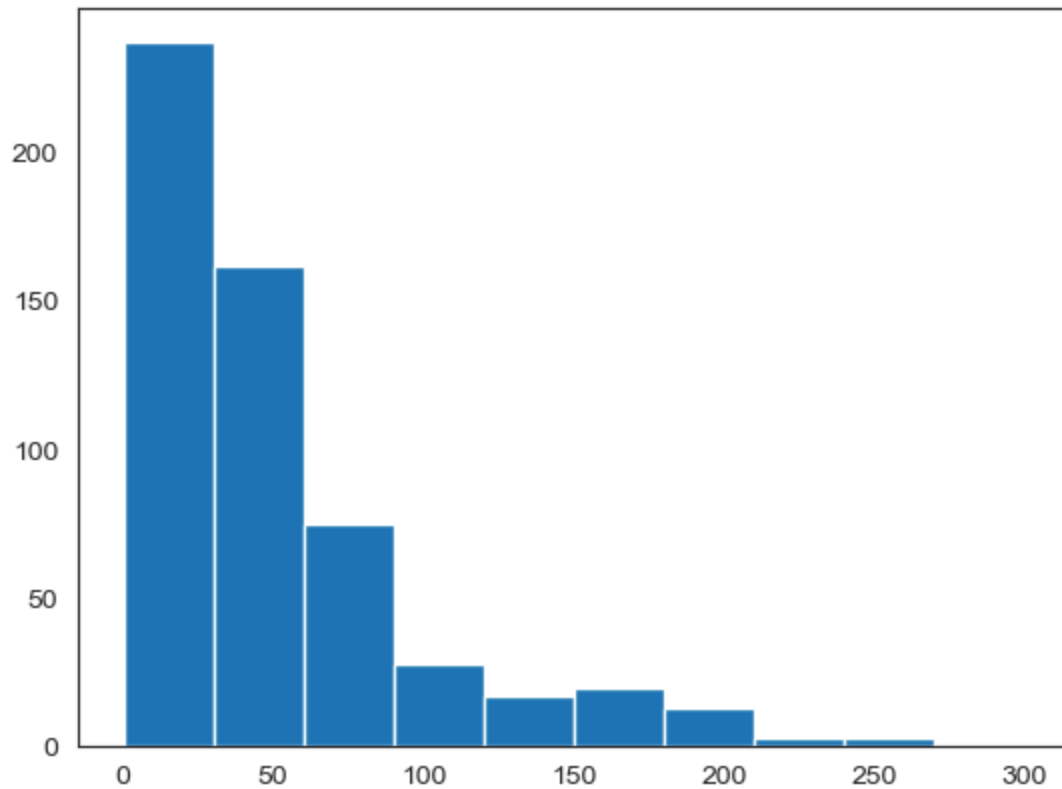


```
In [58]: n1 = plt.hist(movies.AudienceRating, bins = 20)
```

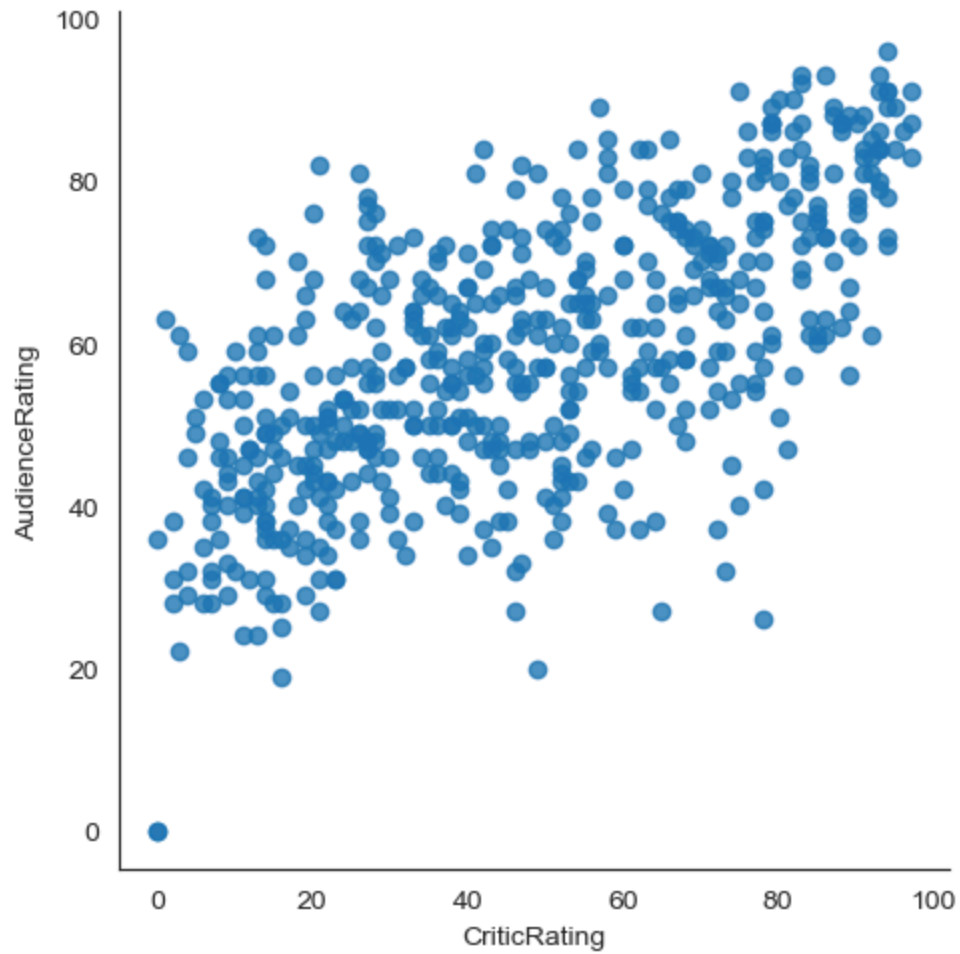


```
In [59]: plt.hist(movies.BudgetMillions)
plt.show
```

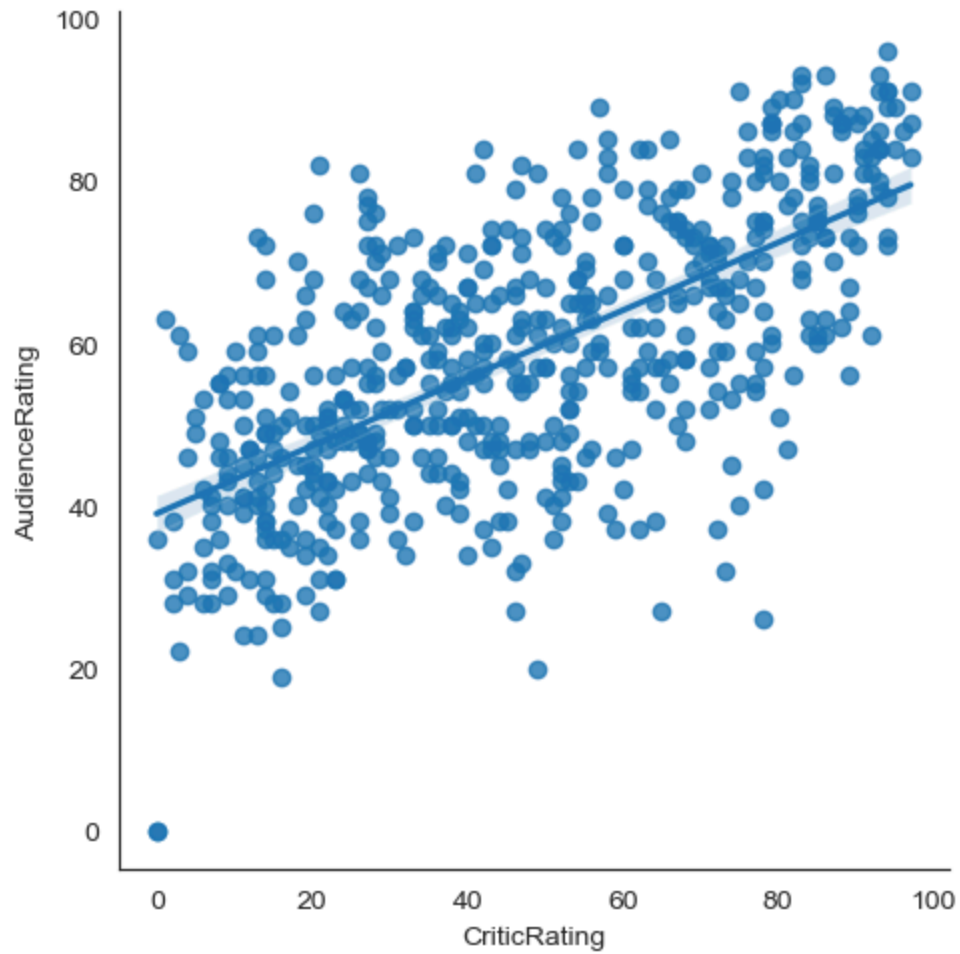
```
Out[59]: <function matplotlib.pyplot.show(close=None, block=None)>
```



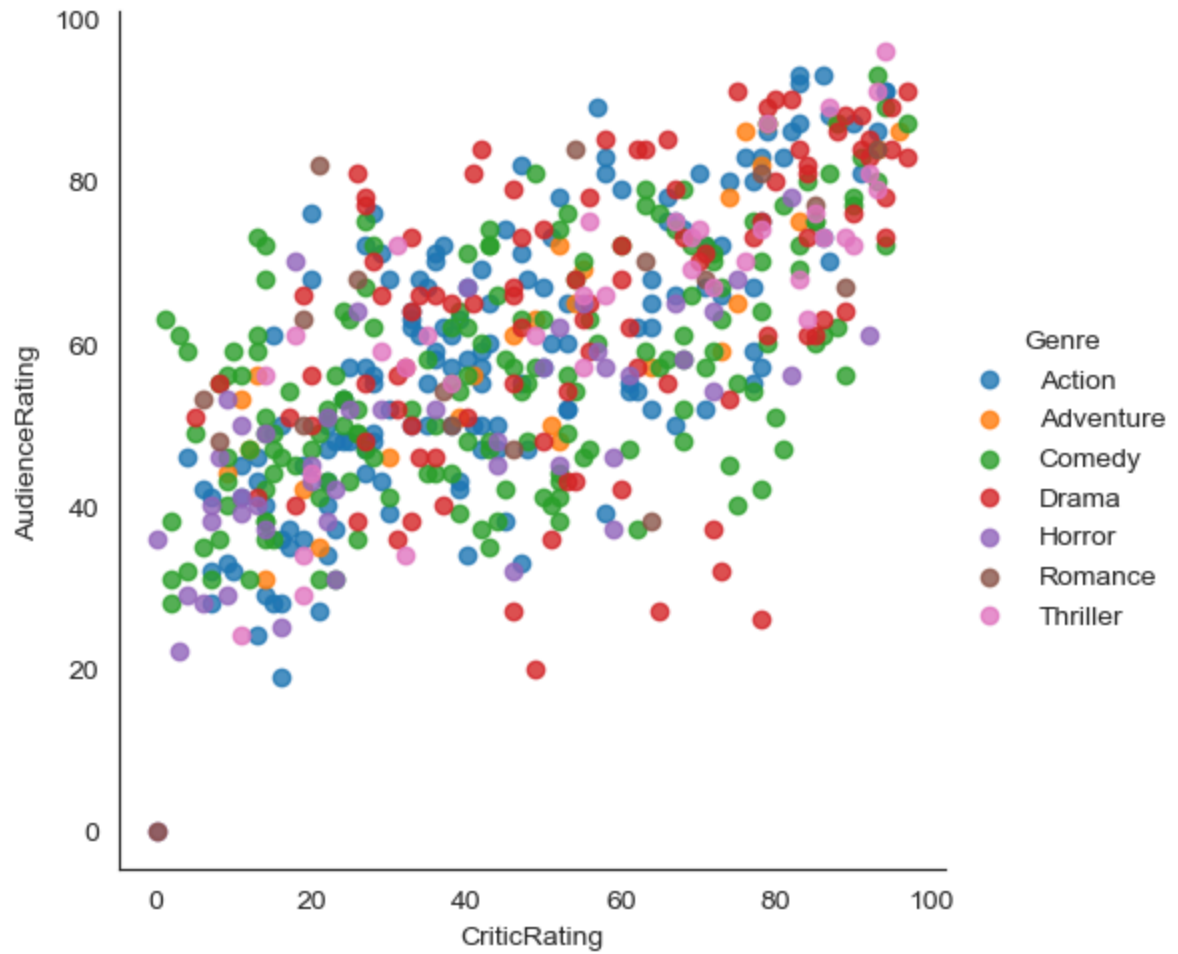
```
In [60]: vis1 = sns.lmplot(data = movies, x = 'CriticRating', y = 'AudienceRating', fit_reg
```



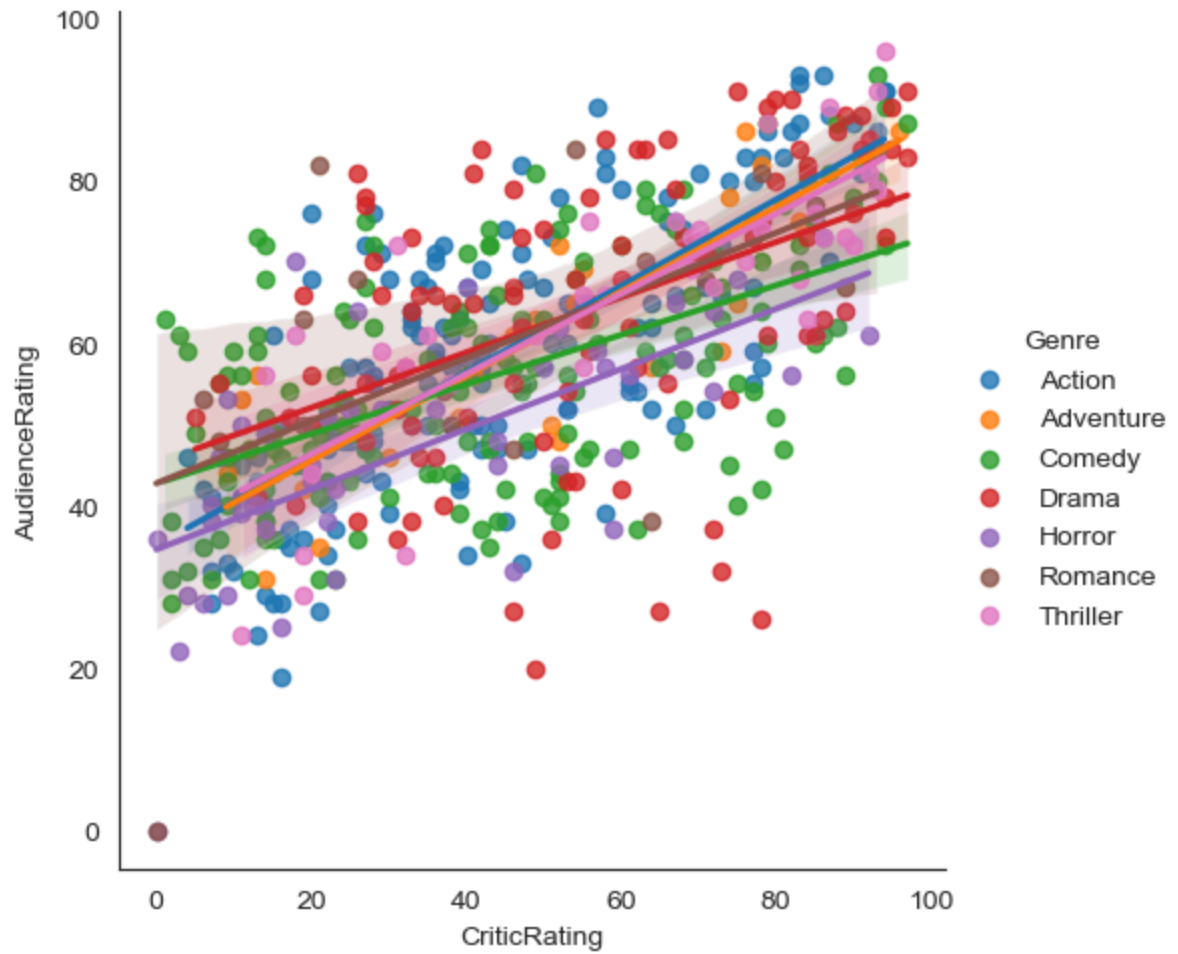
```
In [61]: vis1 = sns.lmplot(data = movies, x = 'CriticRating', y = 'AudienceRating', fit_reg
```



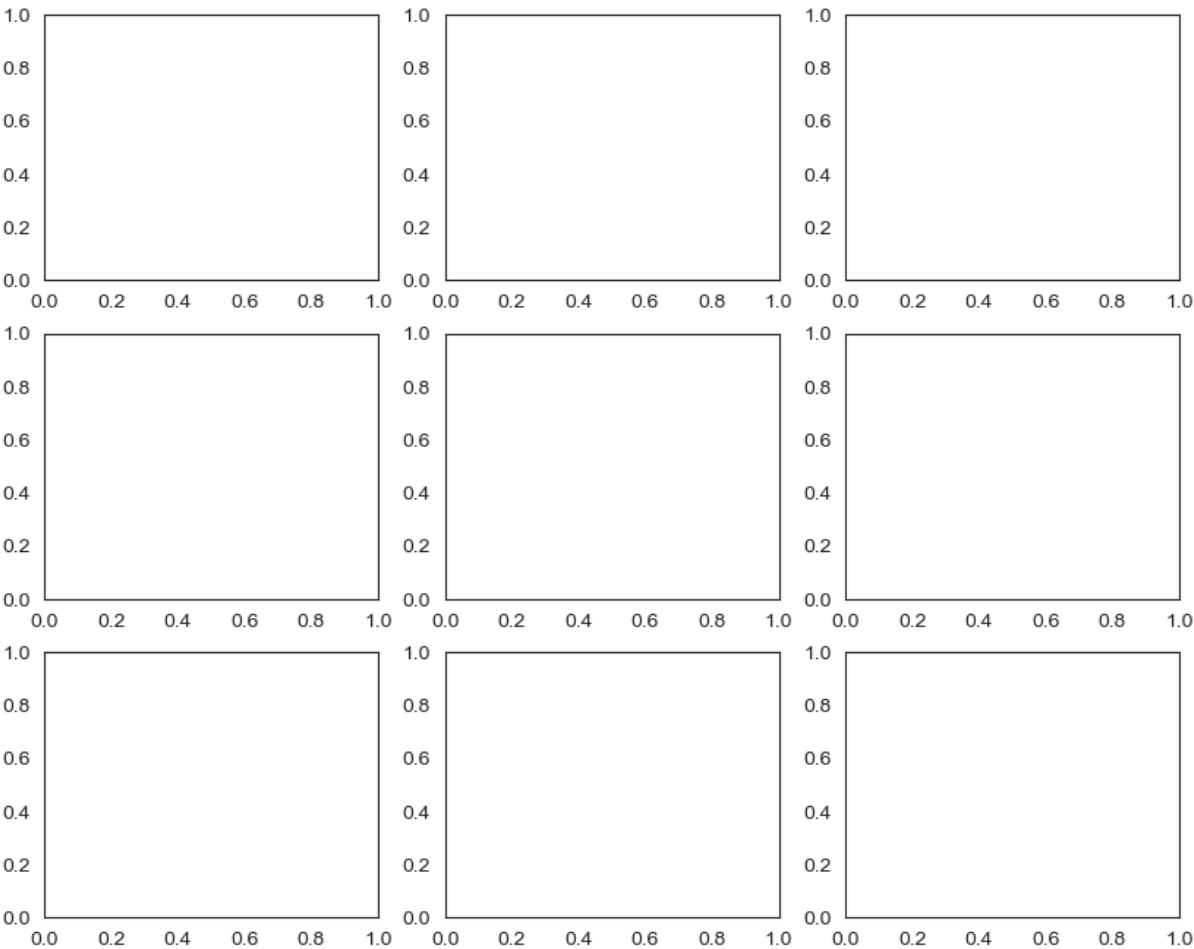
```
In [64]: vis1 = sns.lmplot(data = movies, x = 'CriticRating', y = 'AudienceRating', fit_reg
```

```
In [65]: vis1 = sns.lmplot(data = movies, x = 'CriticRating', y = 'AudienceRating', fit_reg
```



```
In [66]: # subplots
#ax = plt.subplots(1,2, figsize = (3,3,))
ax = plt.subplots(3,3, figsize = (10,8))
```



```
In [ ]:
```