

```
In [1]: import pandas as pd
import os
```

```
In [2]: movies=pd.read_csv(r'D:\Naresh IT foundation\Movie-Rating.csv')
```

```
In [3]: os.getcwd()
```

```
Out[3]: 'C:\\Users\\Rupesh Gupta'
```

```
In [4]: movies
```

```
Out[4]:
```

	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 [5]: type(movies)
```

```
Out[5]: pandas.core.frame.DataFrame
```

```
In [6]: len(movies)
```

```
Out[6]: 559
```

```
In [7]: import numpy
print(numpy.__version__)
```

1.26.4

```
In [8]: import pandas
print(pandas.__version__)
```

2.2.2

```
In [9]: movies.columns
```

```
Out[9]: Index(['Film', 'Genre', 'Rotten Tomatoes Ratings %', 'Audience Ratings %',
   'Budget (million $)', 'Year of release'],
   dtype='object')
```

```
In [10]: 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 [11]: movies.shape#(no of row and no of coloums)
```

```
Out[11]: (559, 6)
```

```
In [12]: movies.head()
```

	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 [13]: movies.tail()
```

Out[13]:

	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 [14]: movies.columns

```
Out[14]: Index(['Film', 'Genre', 'Rotten Tomatoes Ratings %', 'Audience Ratings %',
       'Budget (million $)', 'Year of release'],
       dtype='object')
```

In [15]: movies.columns=['Film', 'Genre', 'CriticRating', 'AudienceRating','BudgetMillions',

In [16]: movies.head(1)*# Removed spaces & % removed noise characters*

```
Out[16]:      Film   Genre  CriticRating  AudienceRating  BudgetMillions  Year
0  (500) Days of Summer  Comedy          87             81            8  2009
```

In [17]: movies.shape

Out[17]: (559, 6)

```
In [18]: movies.describe()
# if you look at the year the data type is int but when you look at the mean value
# we have to change to category type
# also from object datatype we will convert to category datatypes
```

```
Out[18]:      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 [19]: movies.Film=movies.Film.astype('category')
```

```
In [20]: 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 [21]: movies['Film']
#movies['Audience Ratings %']
```

```
Out[21]: 0      (500) Days of Summer
          1      10,000 B.C.
          2      12 Rounds
          3      127 Hours
          4      17 Again
          ...
          554     Your Highness
          555     Youth in Revolt
          556     Zodiac
          557     Zombieland
          558     Zookeeper
Name: Film, Length: 559, dtype: category
Categories (559, object): ['(500) Days of Summer ', '10,000 B.C.', '12 Rounds ',
 '127 Hours', ..., 'Youth in Revolt', 'Zodiac', 'Zombieland ', 'Zookeeper']
```

```
In [22]: movies.Film
```

```
Out[22]: 0      (500) Days of Summer
          1      10,000 B.C.
          2      12 Rounds
          3      127 Hours
          4      17 Again
          ...
          554     Your Highness
          555     Youth in Revolt
          556     Zodiac
          557     Zombieland
          558     Zookeeper
Name: Film, Length: 559, dtype: category
Categories (559, object): ['(500) Days of Summer ', '10,000 B.C.', '12 Rounds ',
 '127 Hours', ..., 'Youth in Revolt', 'Zodiac', 'Zombieland ', 'Zookeeper']
```

```
In [23]: movies.Film = movies.Film.astype('category')
```

In [24]: `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 [25]: `movies.head()`

	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
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 [26]: `movies.Genre=movies.Genre.astype('category')`
`movies.Year=movies.Year.astype('category')`

In [27]: `movies.Genre`

```
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 [28]: `movies.Year# is it real no. year you can take average,min,max but out come have no`

```
Out[28]: 0      2009
         1      2008
         2      2009
         3      2010
         4      2009
         ...
        554    2011
        555    2009
        556    2007
        557    2009
        558    2011
Name: Year, Length: 559, dtype: category
Categories (5, int64): [2007, 2008, 2009, 2010, 2011]
```

In [29]: `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 [30]: `movies.Genre.cat.categories`

```
Out[30]: Index(['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance',
               'Thriller'],
               dtype='object')
```

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

##now when you see the describt you will get only integer value mean, standard devi

Out[31]:

	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

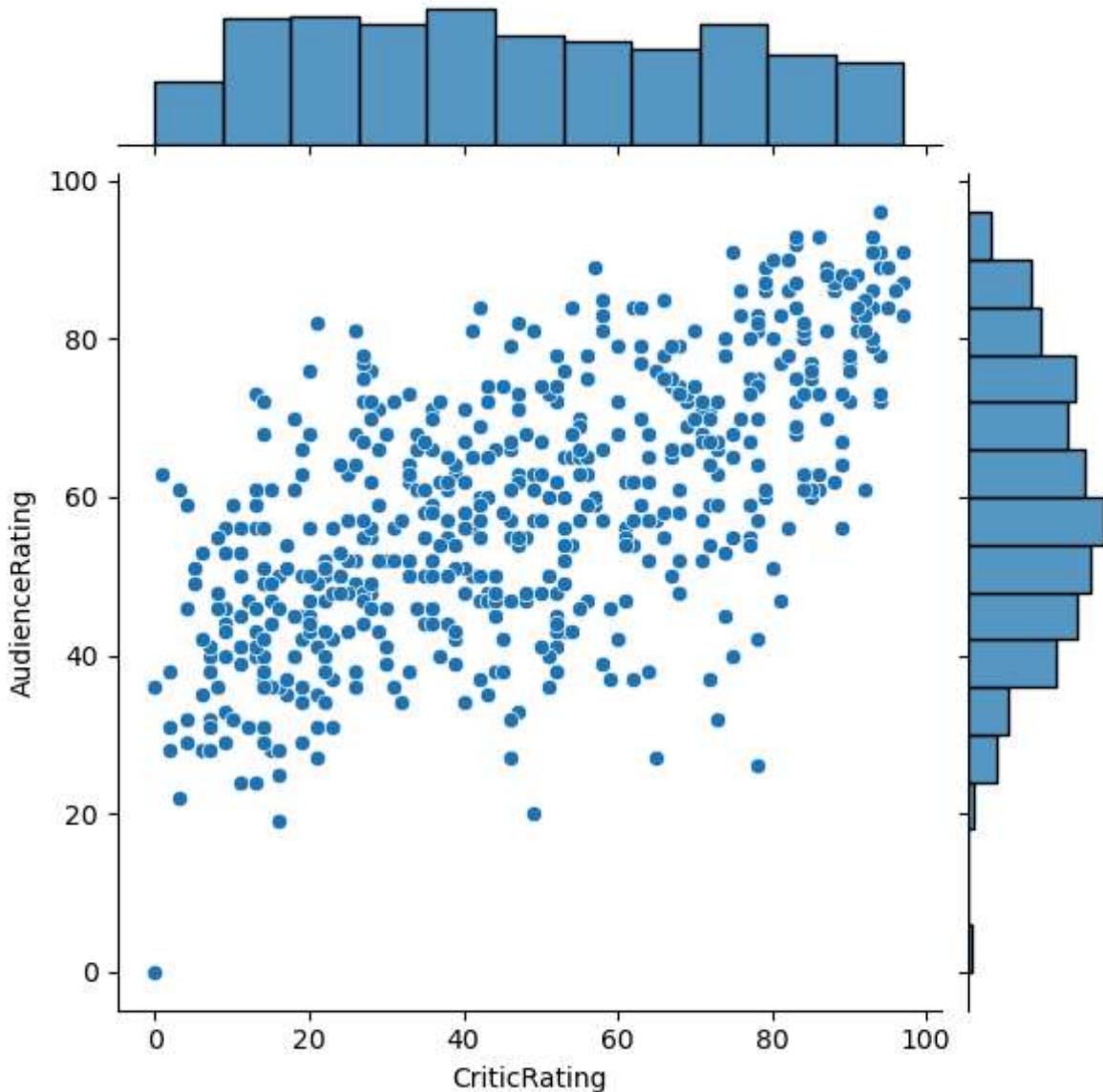
```
In [32]: from matplotlib import pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
In [33]: import warnings
warnings.filterwarnings('ignore')
```

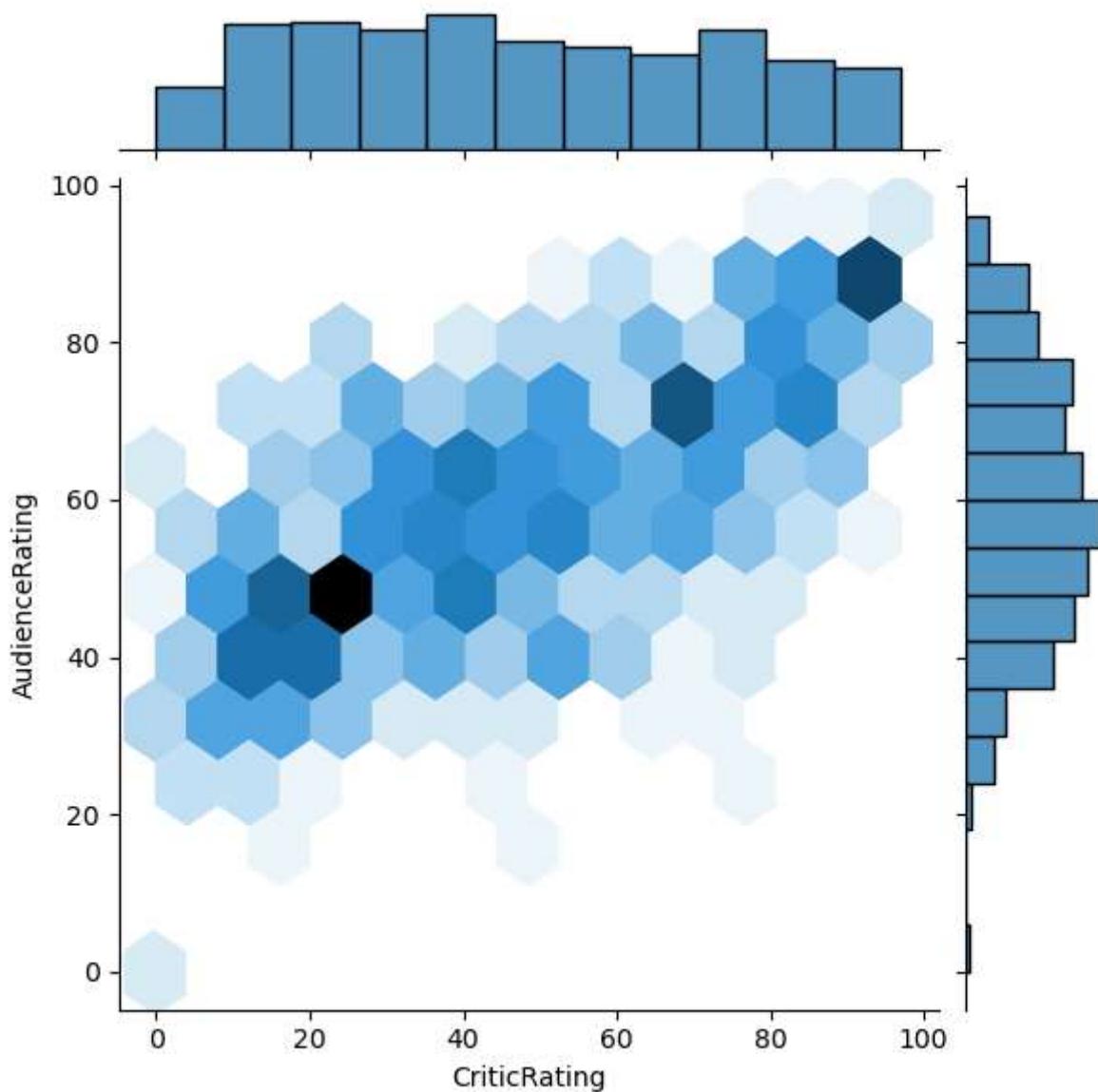
- basically joint plot is a scatter plot & it find the relation b/w audiene & critics
- also if you look up you can find the uniform distribution (critics)and normal distriution (audience)

In [34]:

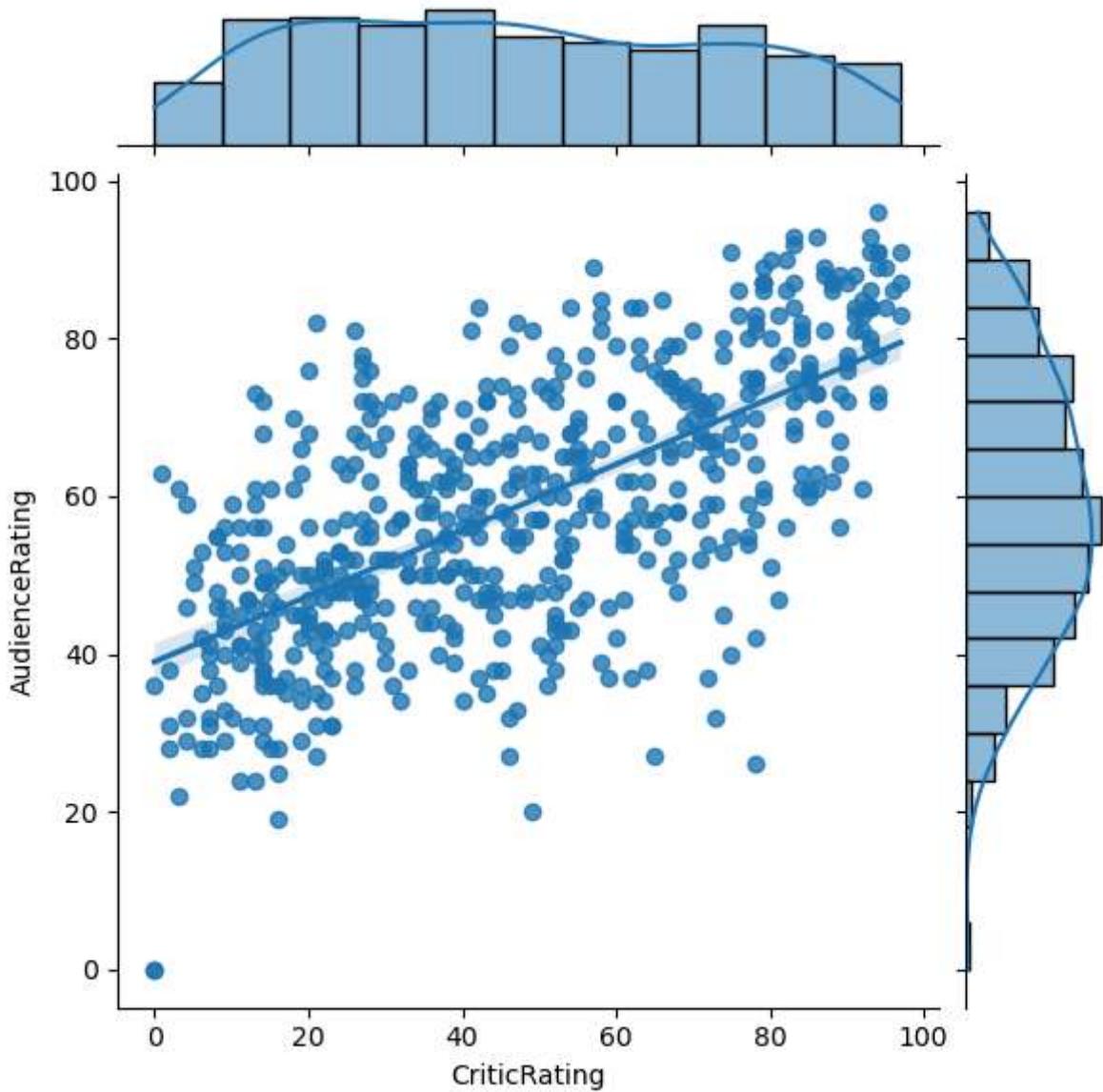
```
j = sns.jointplot( data = movies, x = 'CriticRating', y = 'AudienceRating')
# Audience rating is more dominant than critics rating
# Based on this we find out as most people are most Liklihood to watch audience rat
# Let me explain the excel - if you filter audience rating & critic rating. critic
```



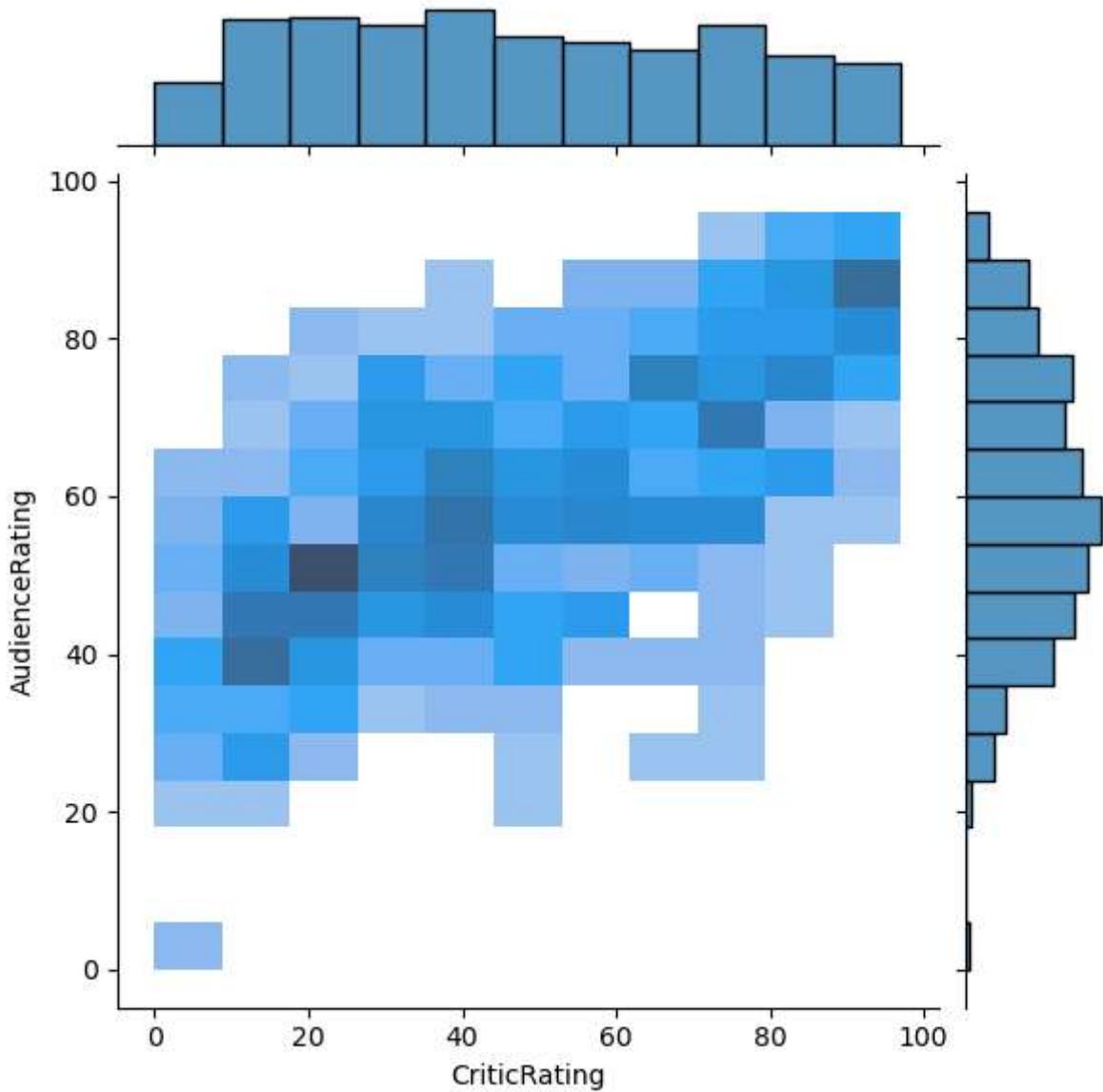
```
In [35]: j=sns.jointplot(data= movies,x='CriticRating',y='AudienceRating',kind='hex')
```



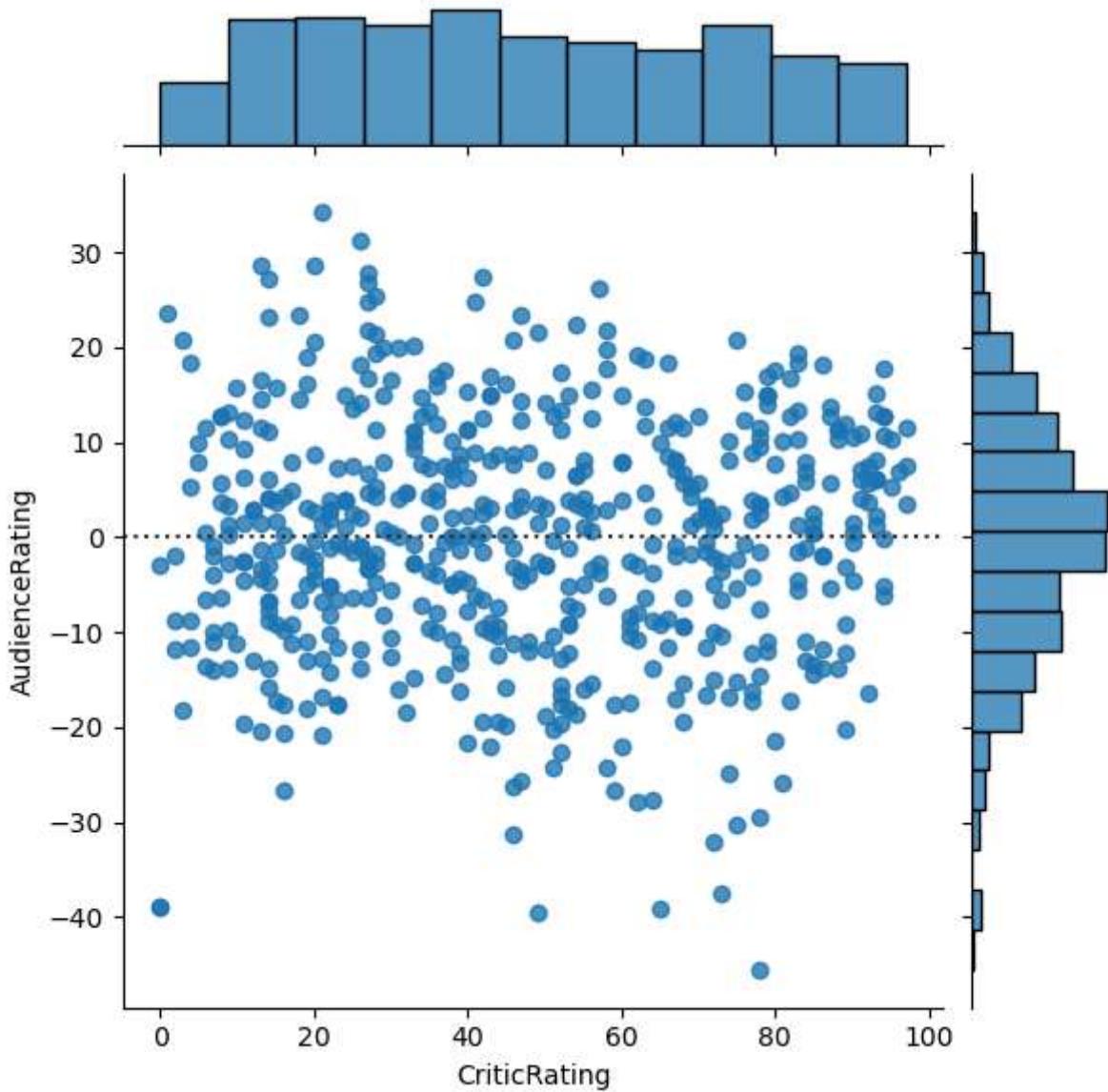
```
In [36]: j = sns.jointplot( data = movies, x = 'CriticRating', y = 'AudienceRating', kind='r'
```



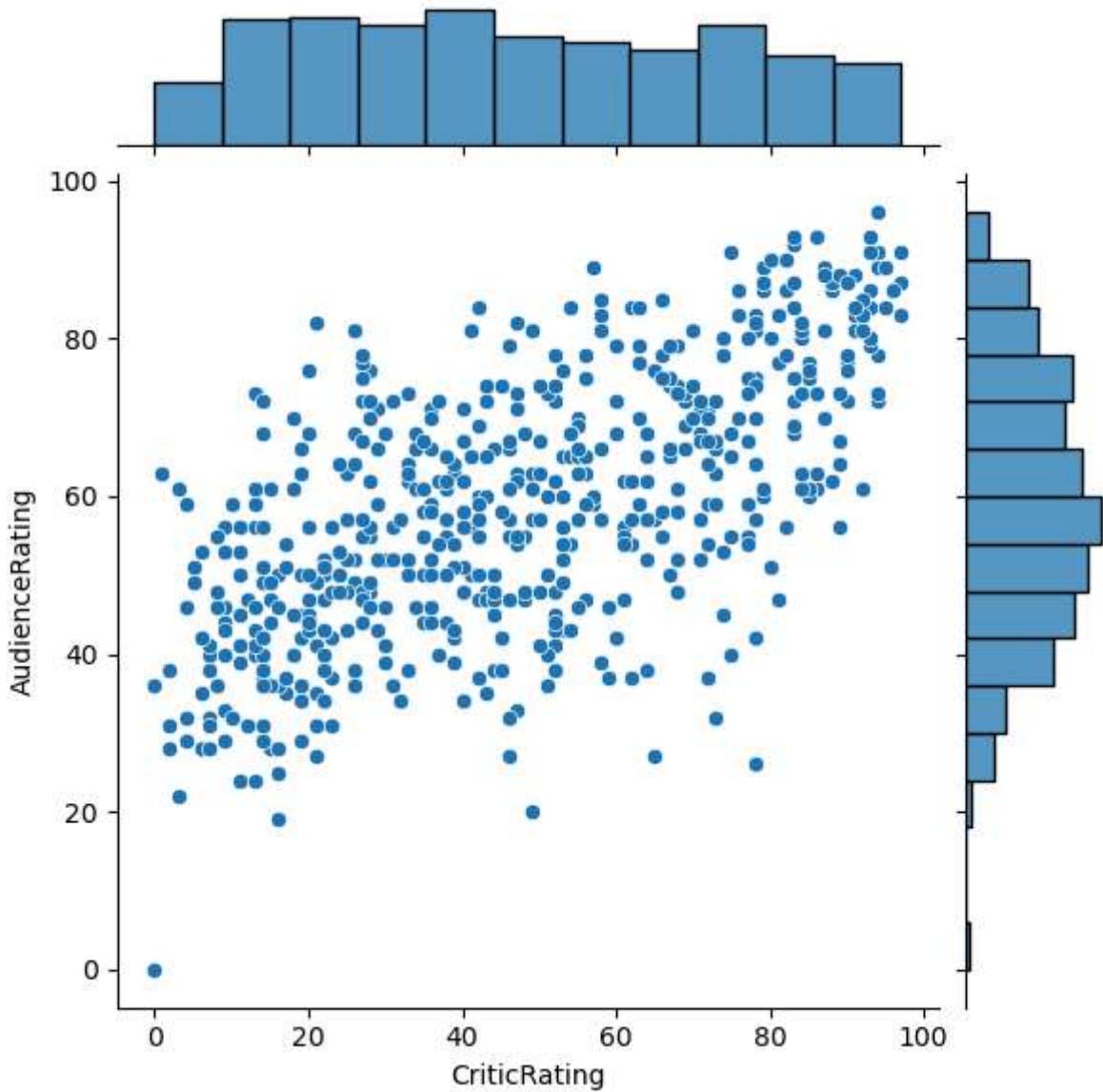
```
In [37]: j=sns.jointplot(data= movies,x='CriticRating',y='AudienceRating',kind='hist')
```



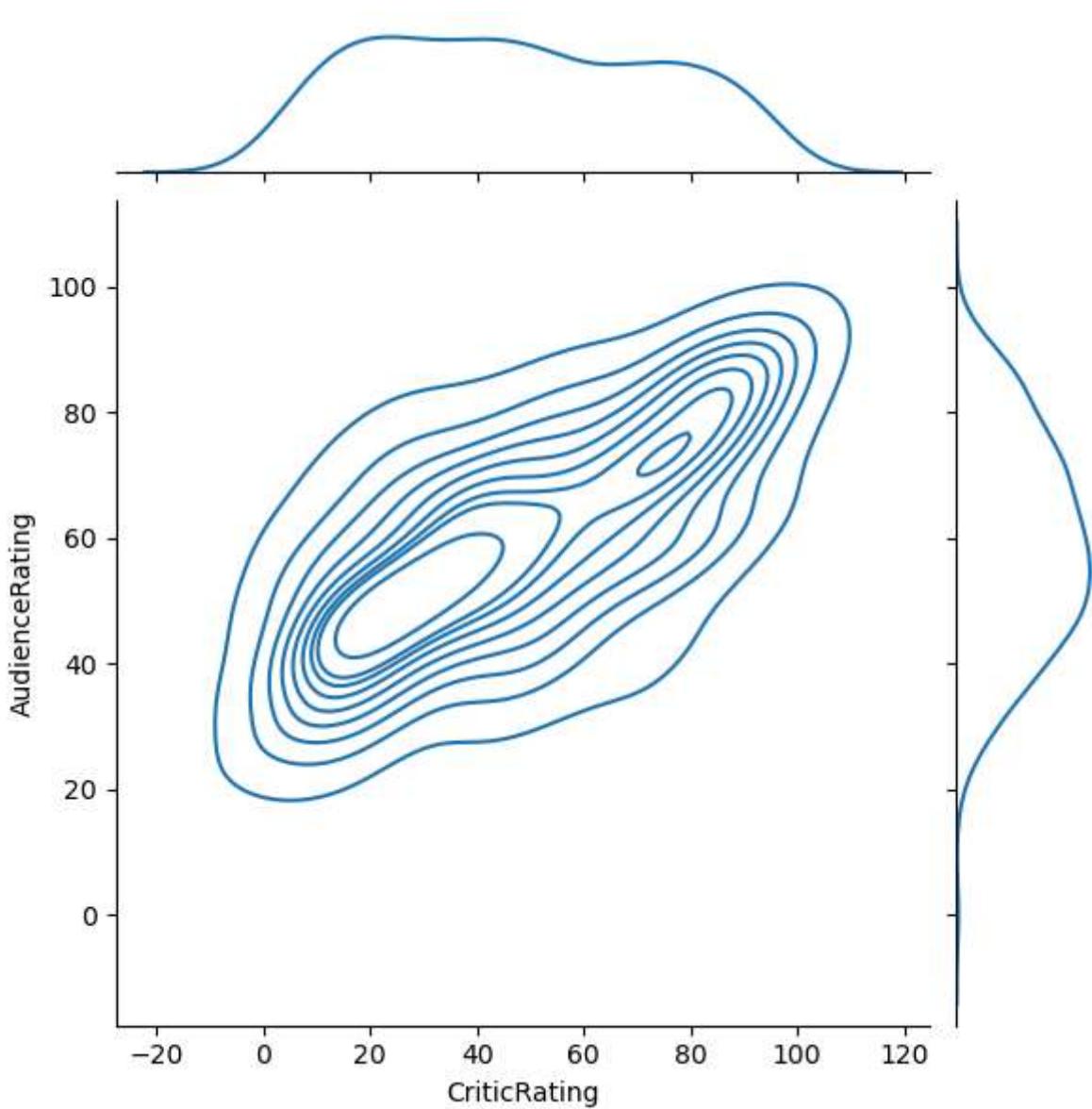
```
In [38]: j=sns.jointplot(data= movies,x='CriticRating',y='AudienceRating',kind='resid')
```



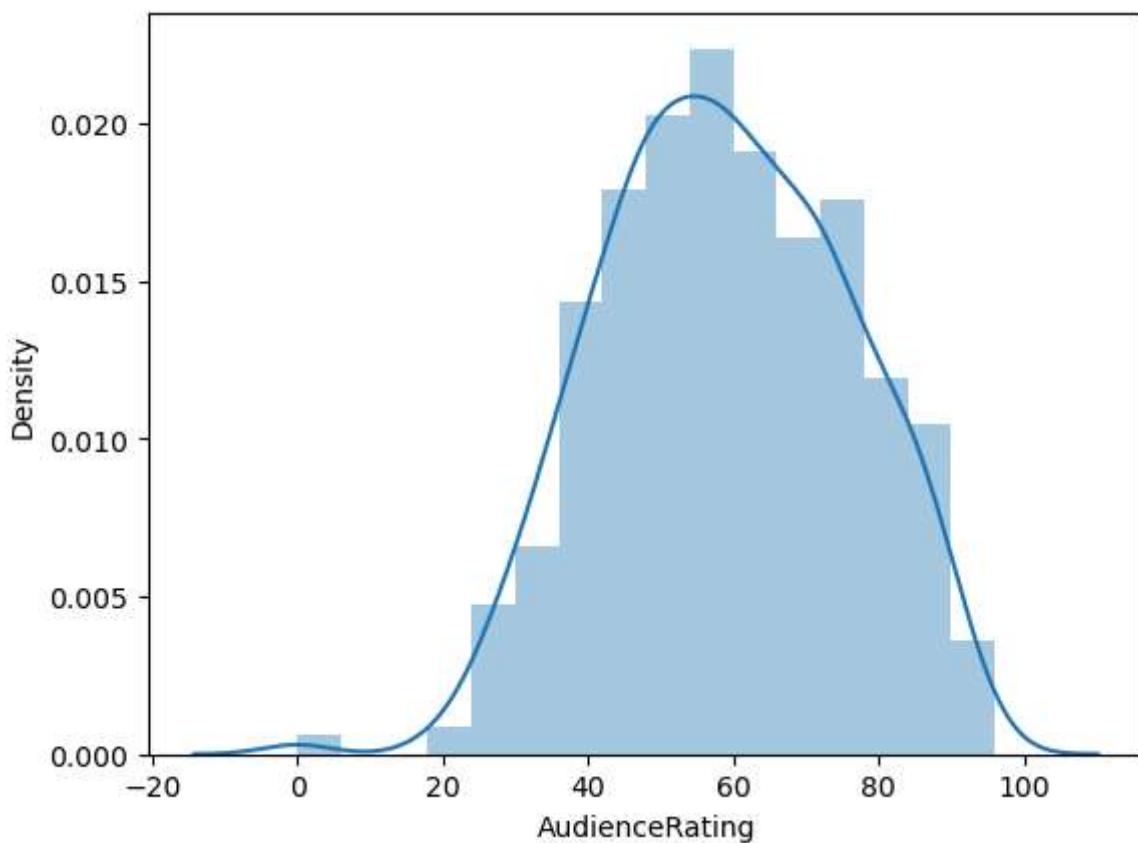
```
In [39]: j=sns.jointplot(data= movies,x='CriticRating',y='AudienceRating',kind='scatter')
```



```
In [40]: j=sns.jointplot(data= movies,x='CriticRating',y='AudienceRating',kind='kde')
```

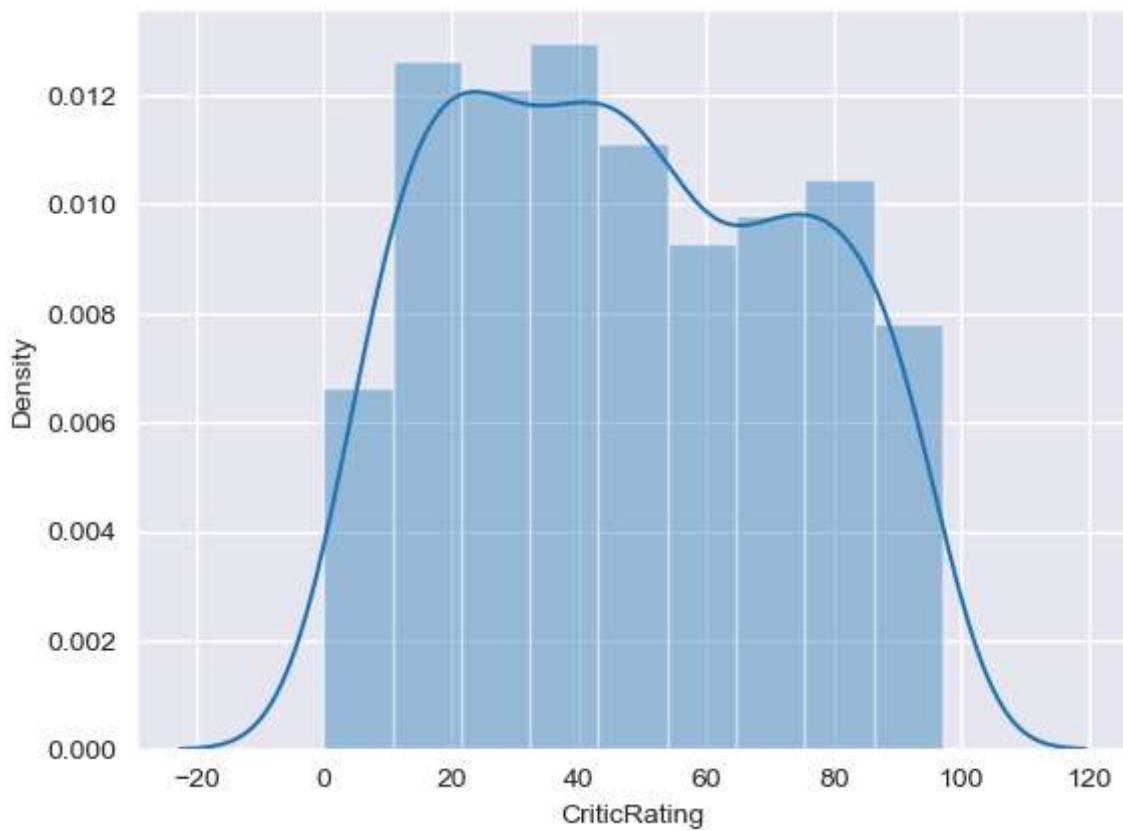


```
In [41]: #Histograms  
# <<< chat1  
m1 = sns.distplot(movies.AudienceRating)  
#y - axis generated by seaborn automatically that is the powefull of seaborn galler
```



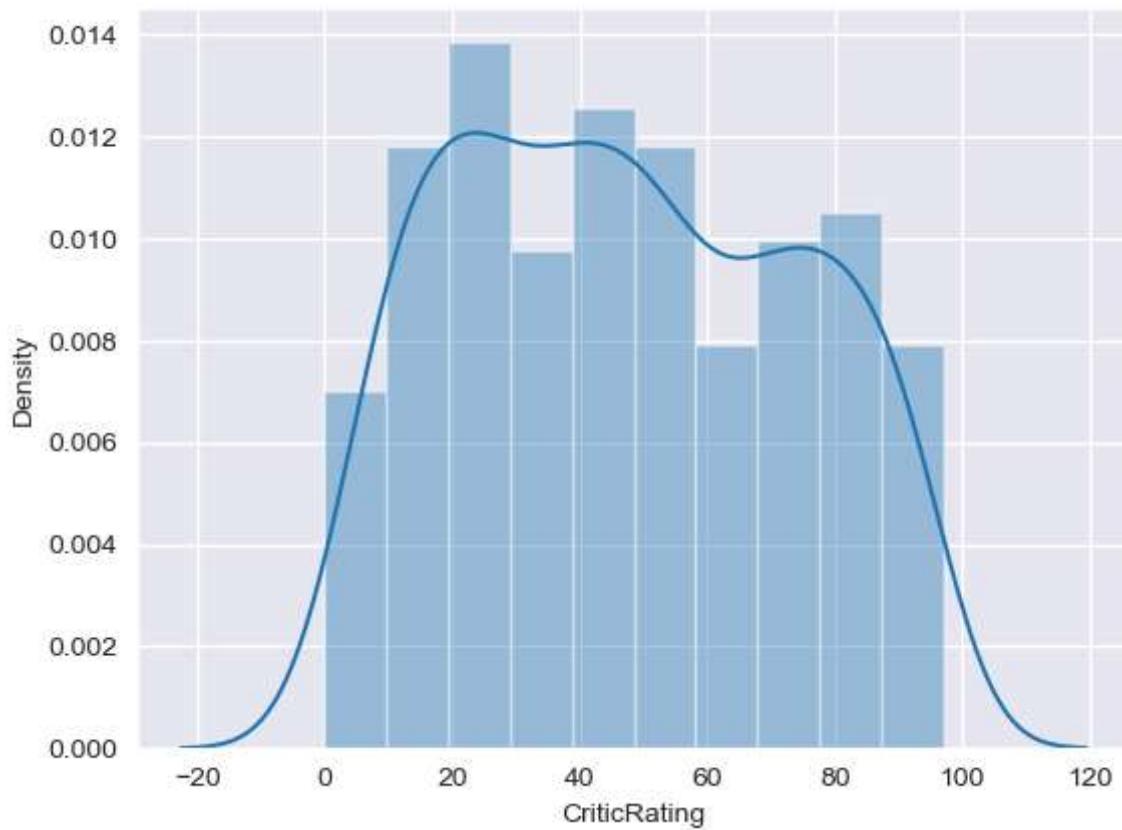
```
In [42]: sns.set_style('darkgrid')
```

```
In [43]: m1=sns.distplot(movies.CriticRating)
```

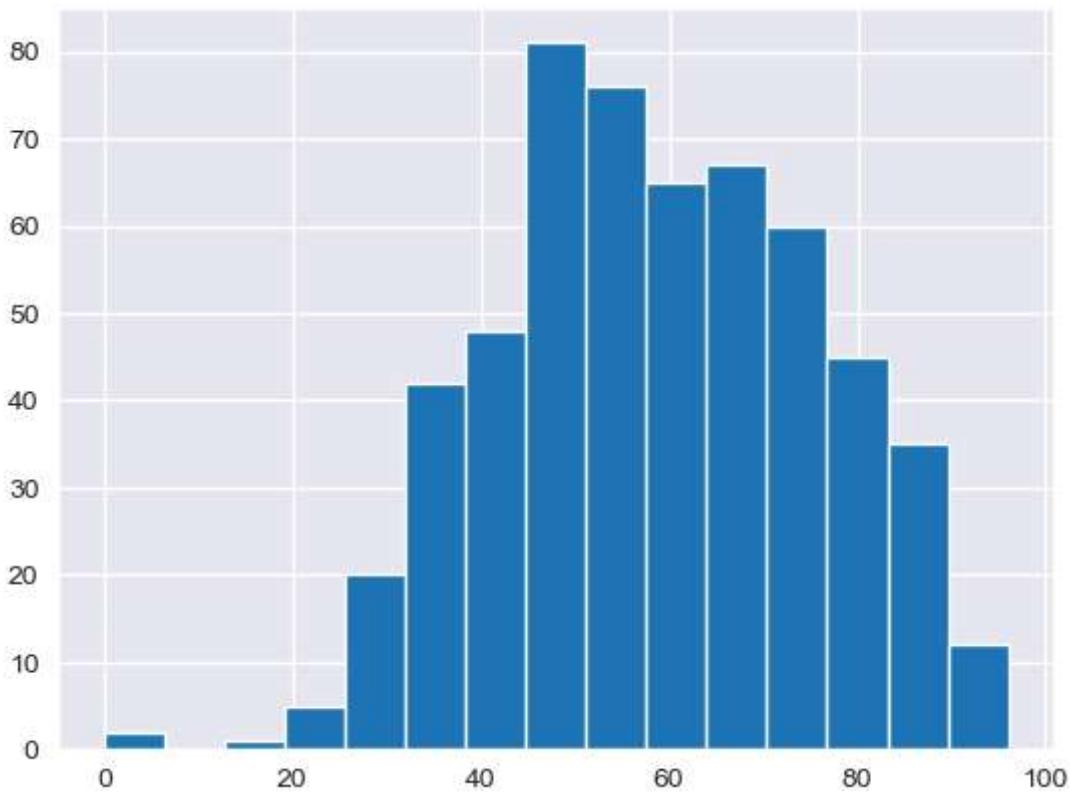


```
In [44]: sns.set_style('darkgrid')
```

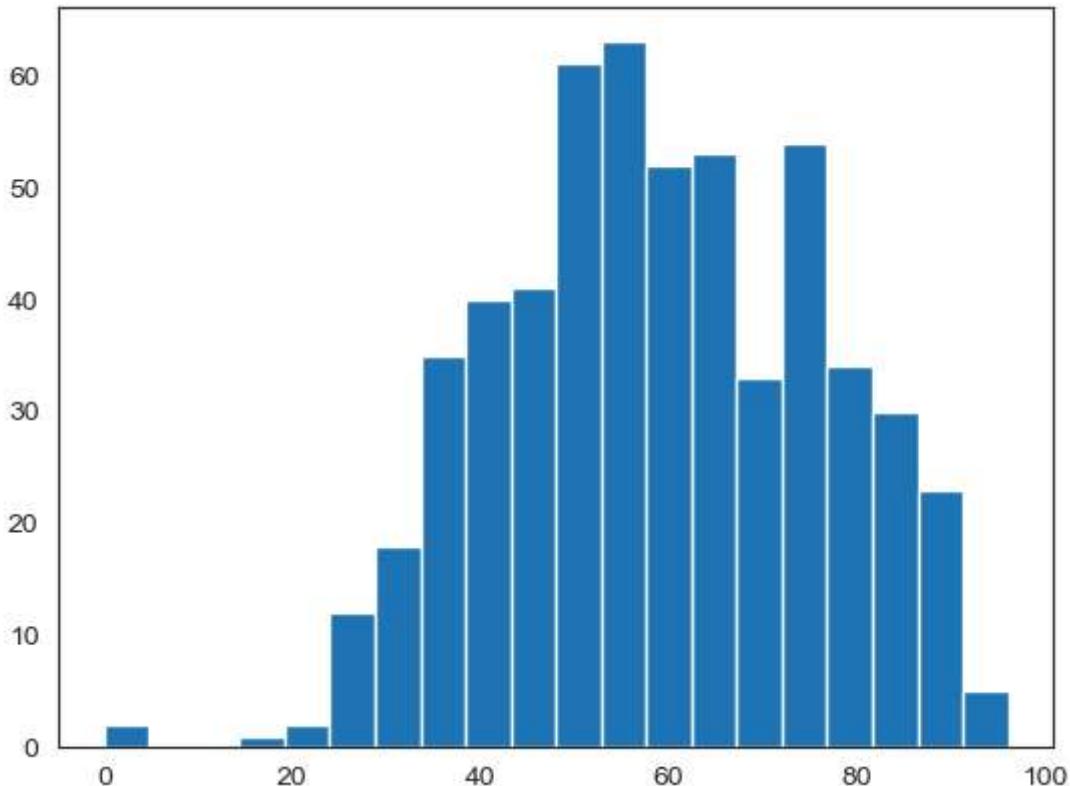
```
In [45]: m2=sns.distplot(movies.CriticRating,bins=10)
```



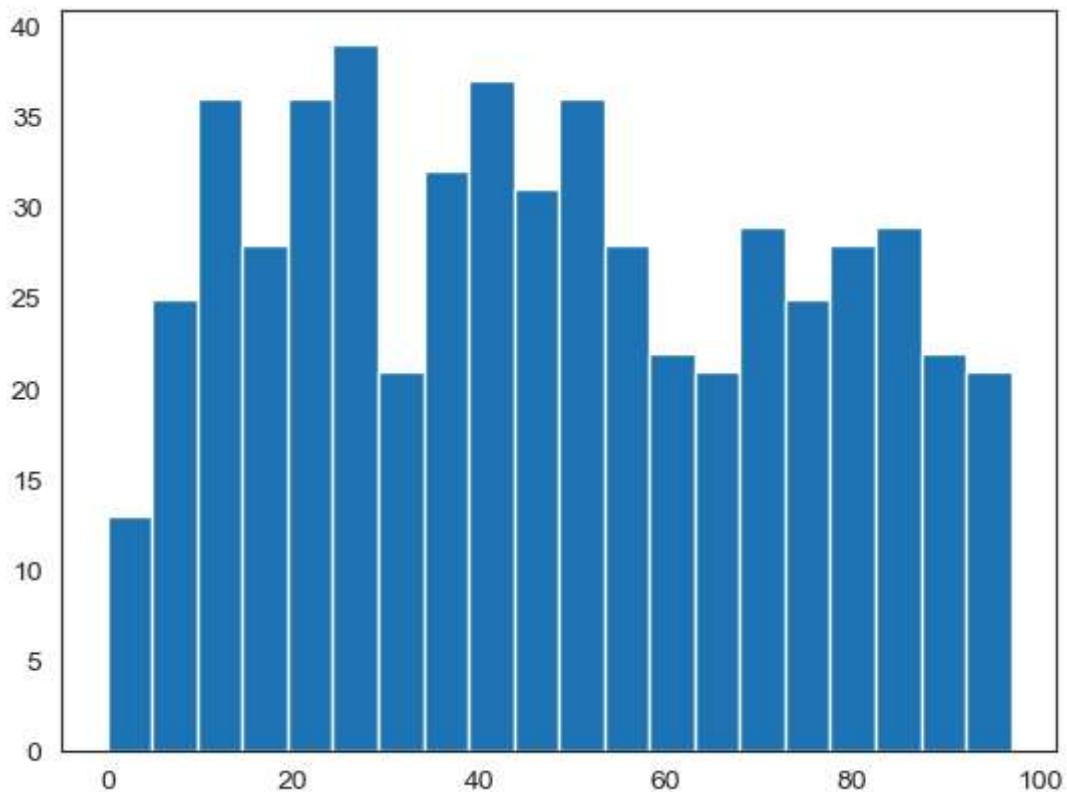
```
In [46]: #sns.set_style('darkgrid')
n1 = plt.hist(movies.AudienceRating, bins=15)
```



```
In [47]: sns.set_style('white') #normal distribution & called as bell curve  
n1 = plt.hist(movies.AudienceRating, bins=20)
```

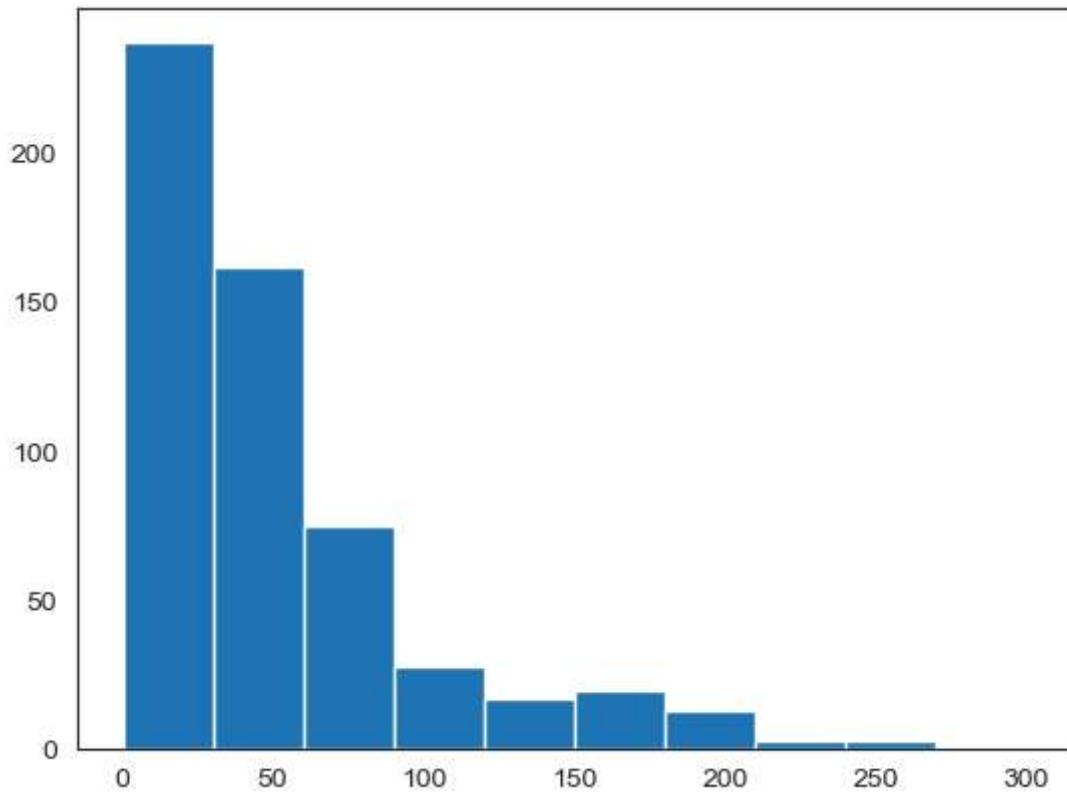


```
In [48]: n1 = plt.hist(movies.CriticRating, bins=20) #uniform distribution
```



```
In [49]: # <<< chat - 2  
# Creating stacked histograms & this is bit tough to understand
```

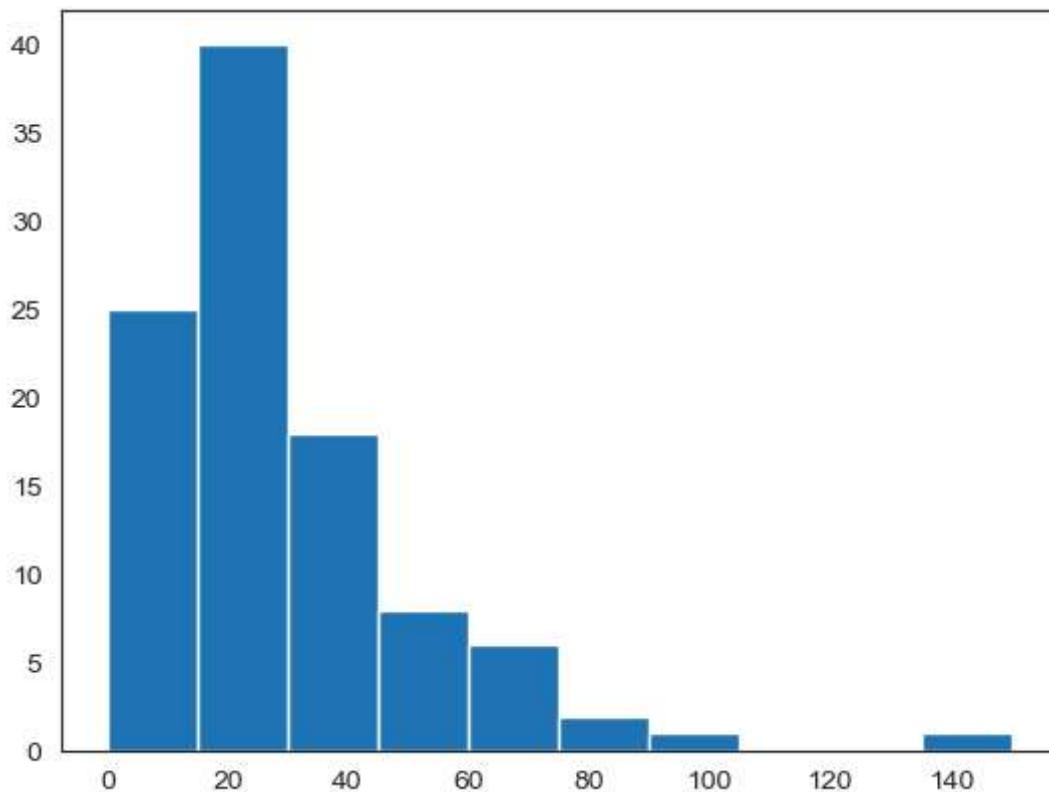
```
In [50]: plt.hist(movies.BudgetMillions)  
plt.show()
```



```
In [51]: movies.columns
```

```
Out[51]: Index(['Film', 'Genre', 'CriticRating', 'AudienceRating', 'BudgetMillions',
       'Year'],
      dtype='object')
```

```
In [52]: plt.hist(movies[movies.Genre == 'Drama'].BudgetMillions)
plt.show()
```



In [53]: `movies.head()`

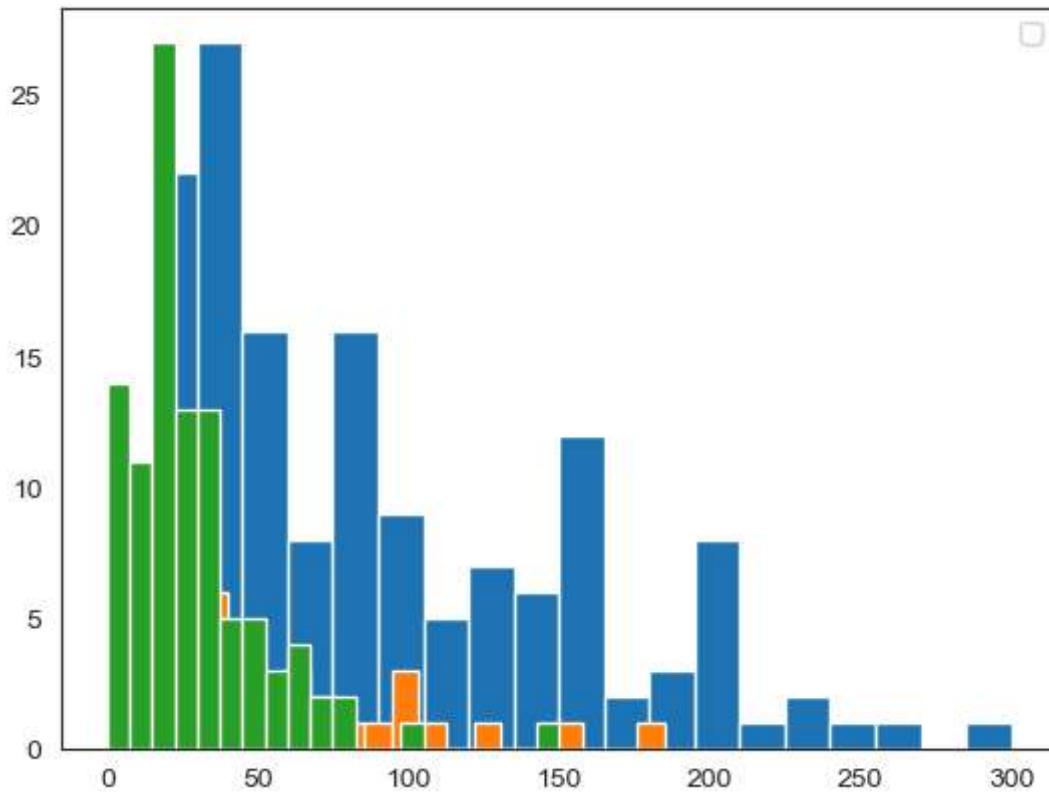
Out[53]:

	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
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

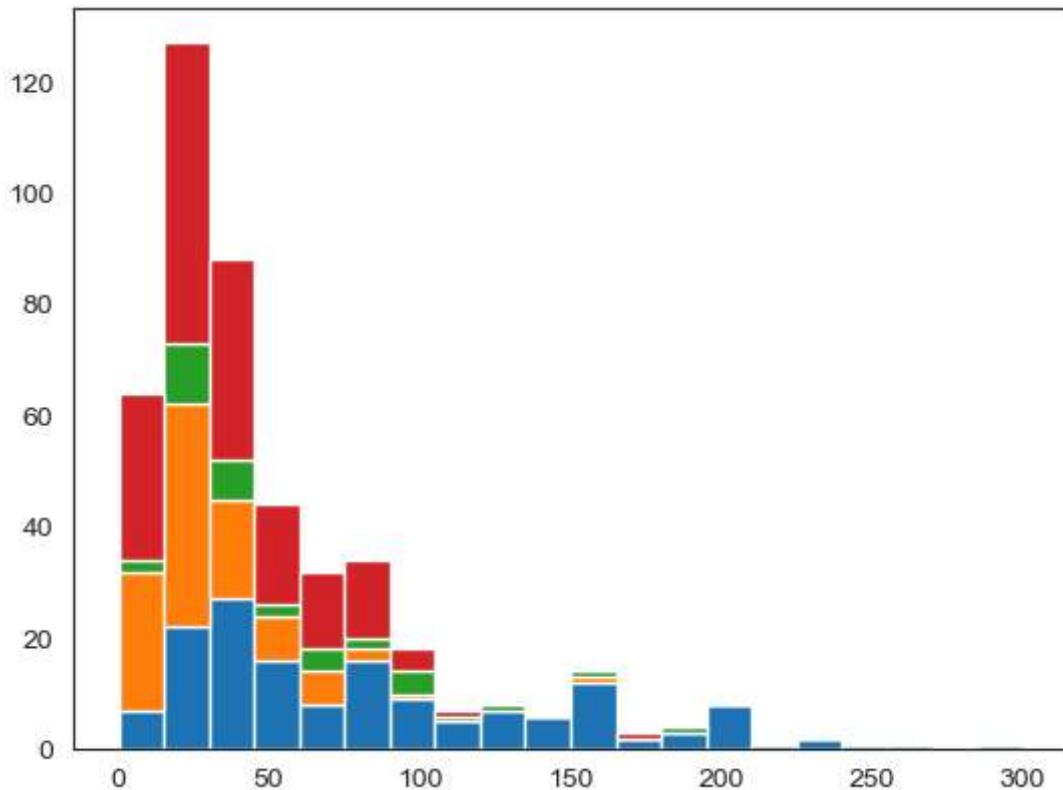
In [54]: `#movies.Genre.unique()`

In [55]: `# Below plots are stacked histogram because overlaped`

```
plt.hist(movies[movies.Genre == 'Action'].BudgetMillions, bins = 20)
plt.hist(movies[movies.Genre == 'Thriller'].BudgetMillions, bins = 20)
plt.hist(movies[movies.Genre == 'Drama'].BudgetMillions, bins = 20)
plt.legend()
plt.show()
```



```
In [56]: plt.hist([movies[movies.Genre == 'Action'].BudgetMillions,\n               movies[movies.Genre == 'Drama'].BudgetMillions, \\\n               movies[movies.Genre == 'Thriller'].BudgetMillions, \\\n               movies[movies.Genre == 'Comedy'].BudgetMillions],\n               bins = 20, stacked = True)\nplt.show()
```



In [57]: movies

Out[57]:

	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
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

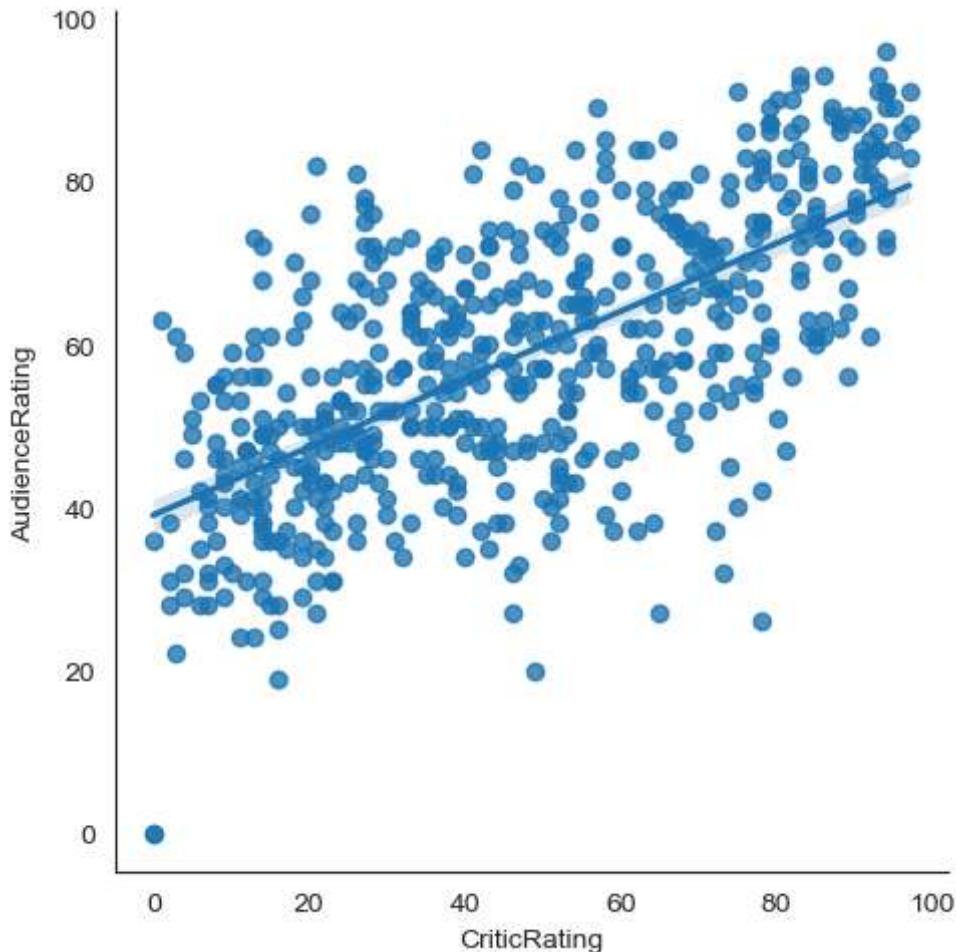
559 rows × 6 columns

In [58]: # if you have 100 categories you cannot copy & paste all the things

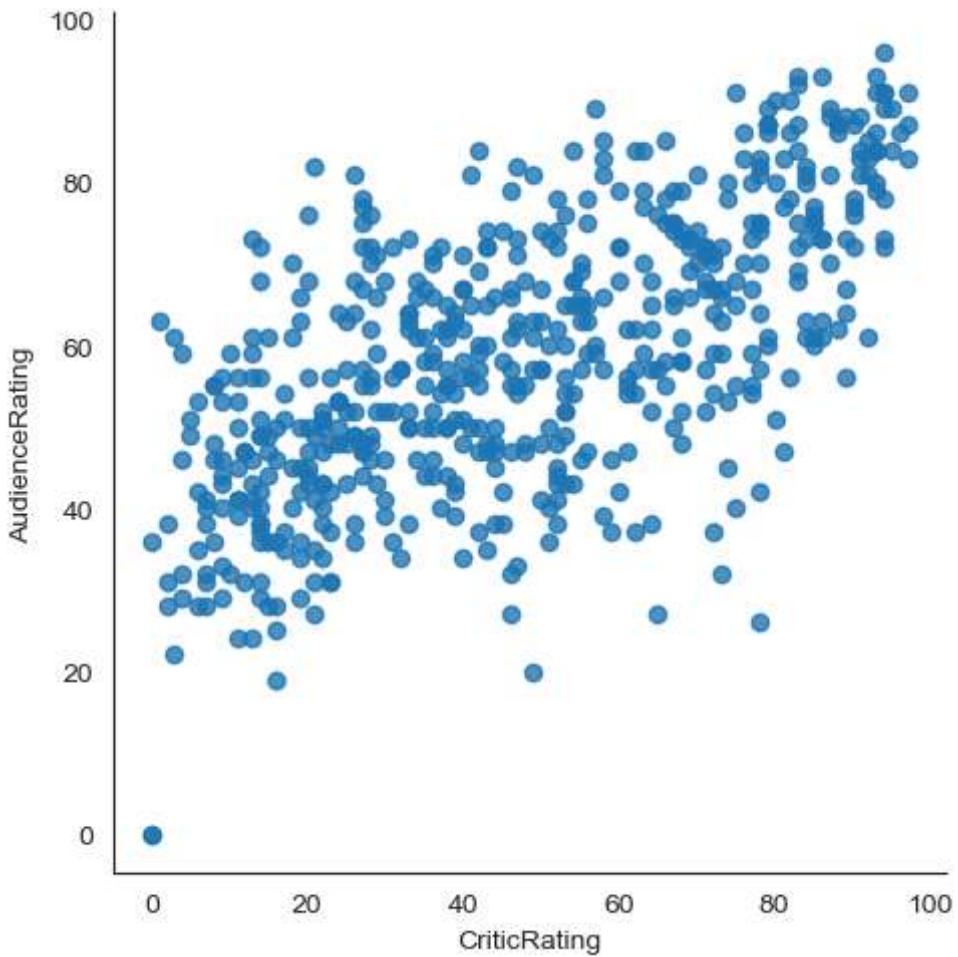
```
for gen in movies.Genre.cat.categories:  
    print(gen)
```

Action
Adventure
Comedy
Drama
Horror
Romance
Thriller

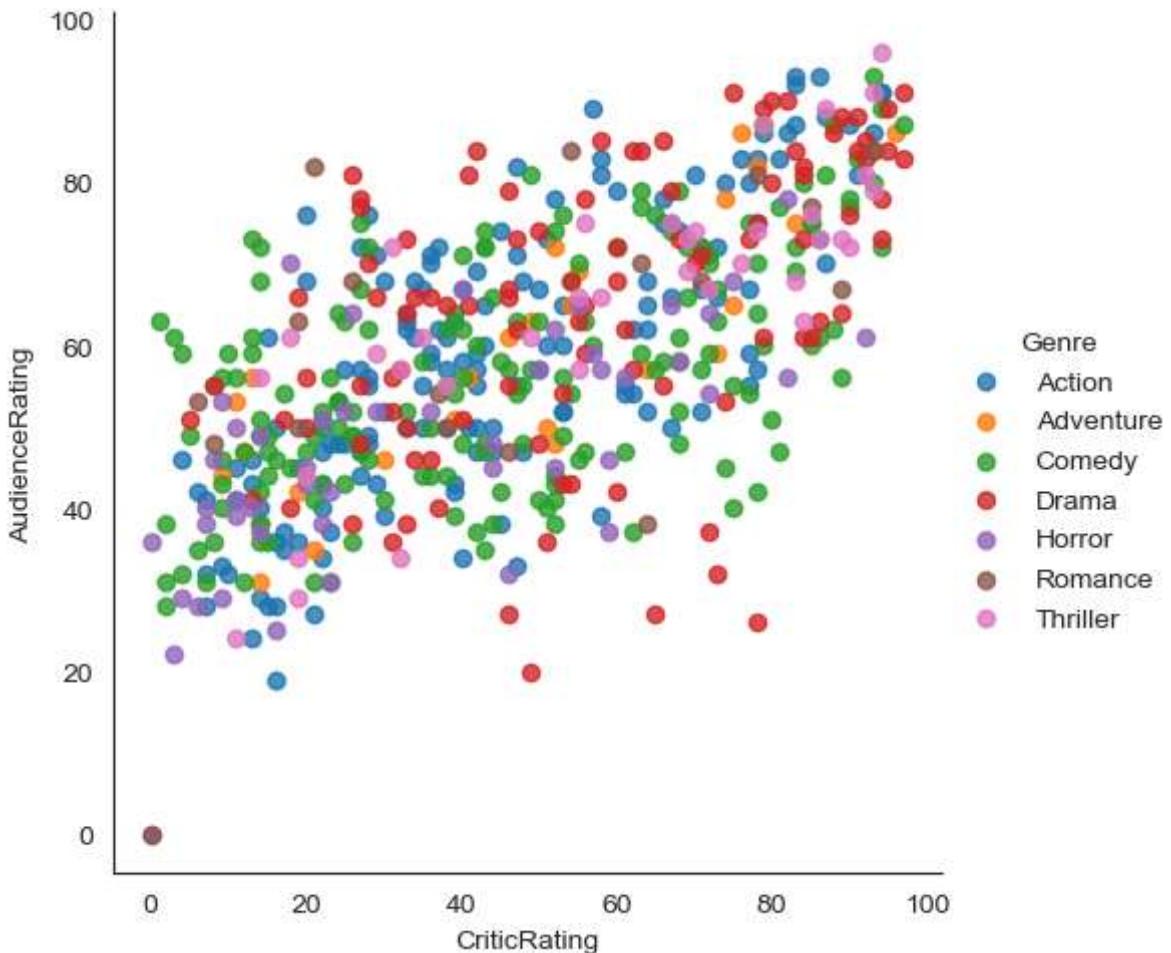
```
In [59]: vis1=sns.lmplot(data=movies, x='CriticRating',y='AudienceRating',  
                      fit_reg=True)
```



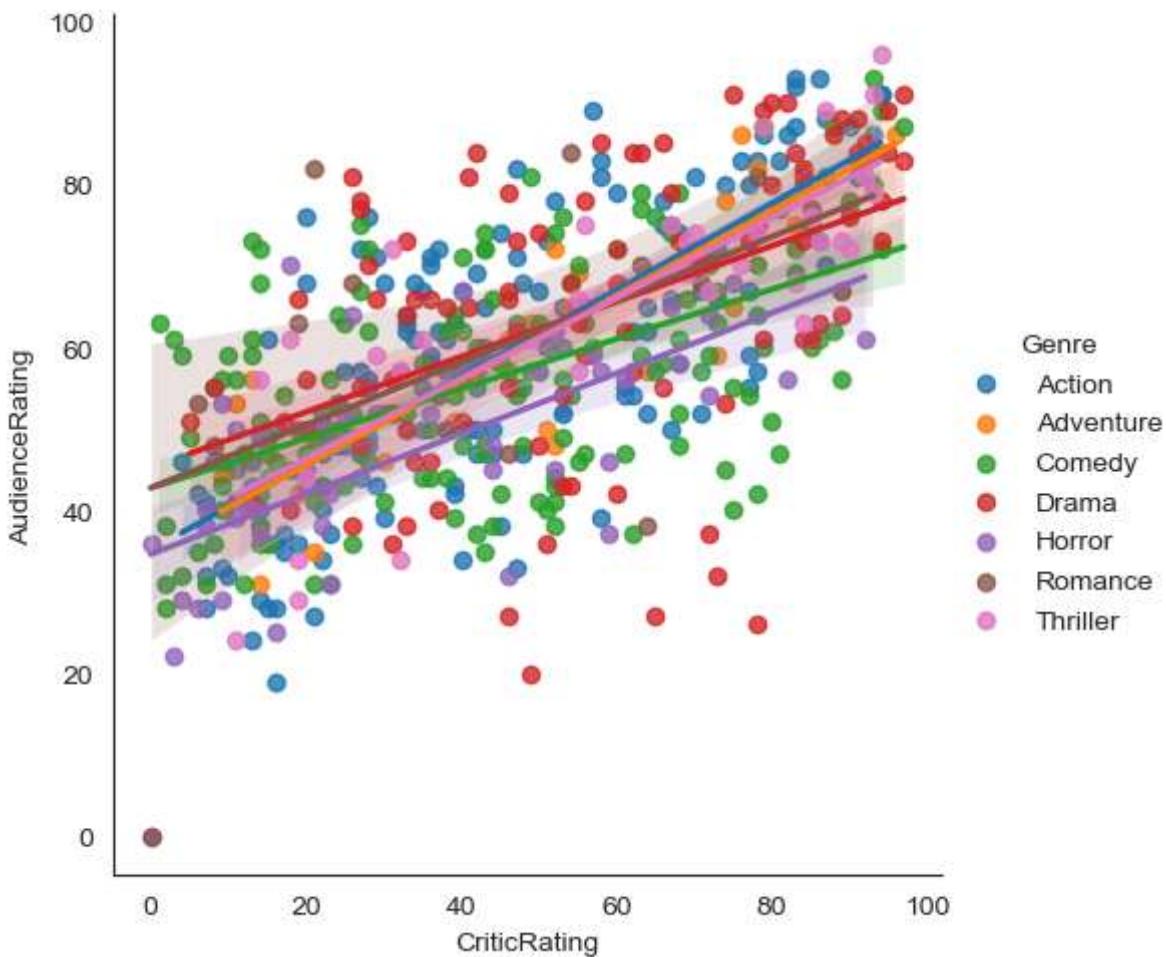
```
In [60]: vis1=sns.lmplot(data=movies, x='CriticRating',y='AudienceRating',fit_reg=False)
```



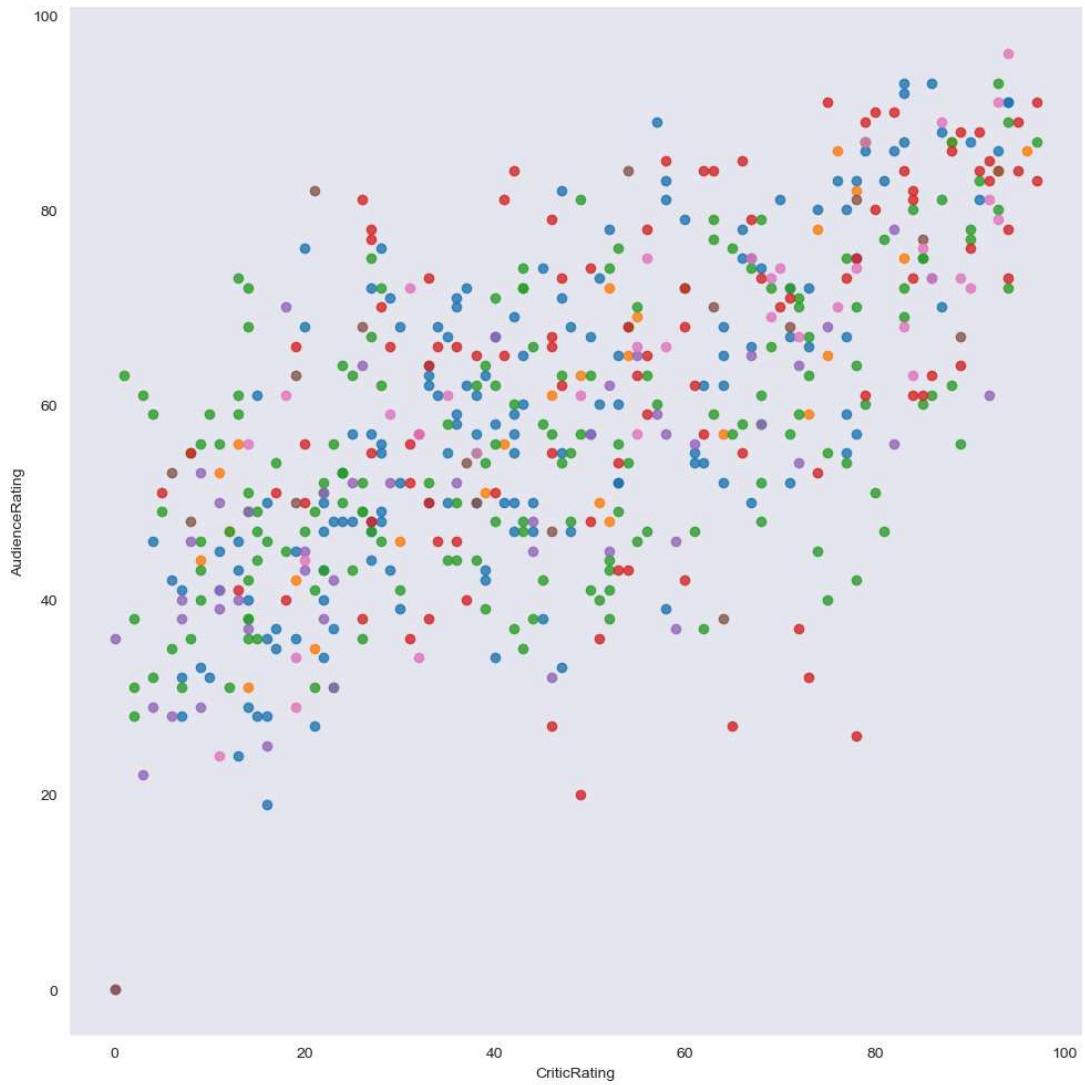
```
In [61]: vis1=sns.lmplot(data=movies, x='CriticRating',y='AudienceRating',fit_reg=False,hue=
```



```
In [62]: vis1=sns.lmplot(data=movies, x='CriticRating',y='AudienceRating',fit_reg=True,hue='
```



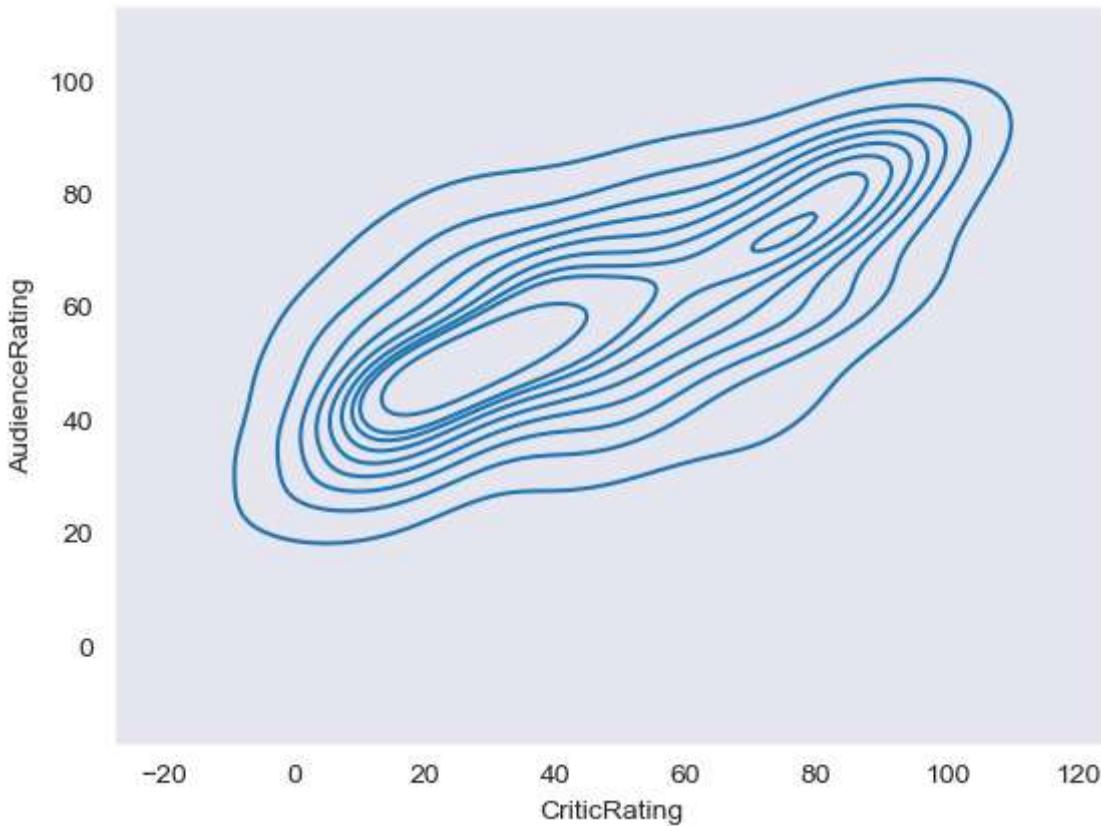
```
In [67]: vis1 = sns.lmplot(data=movies, x='CriticRating', y='AudienceRating',\n                      fit_reg=False, hue = 'Genre', height =10, aspect=1)\nsns.set_style('dark')
```



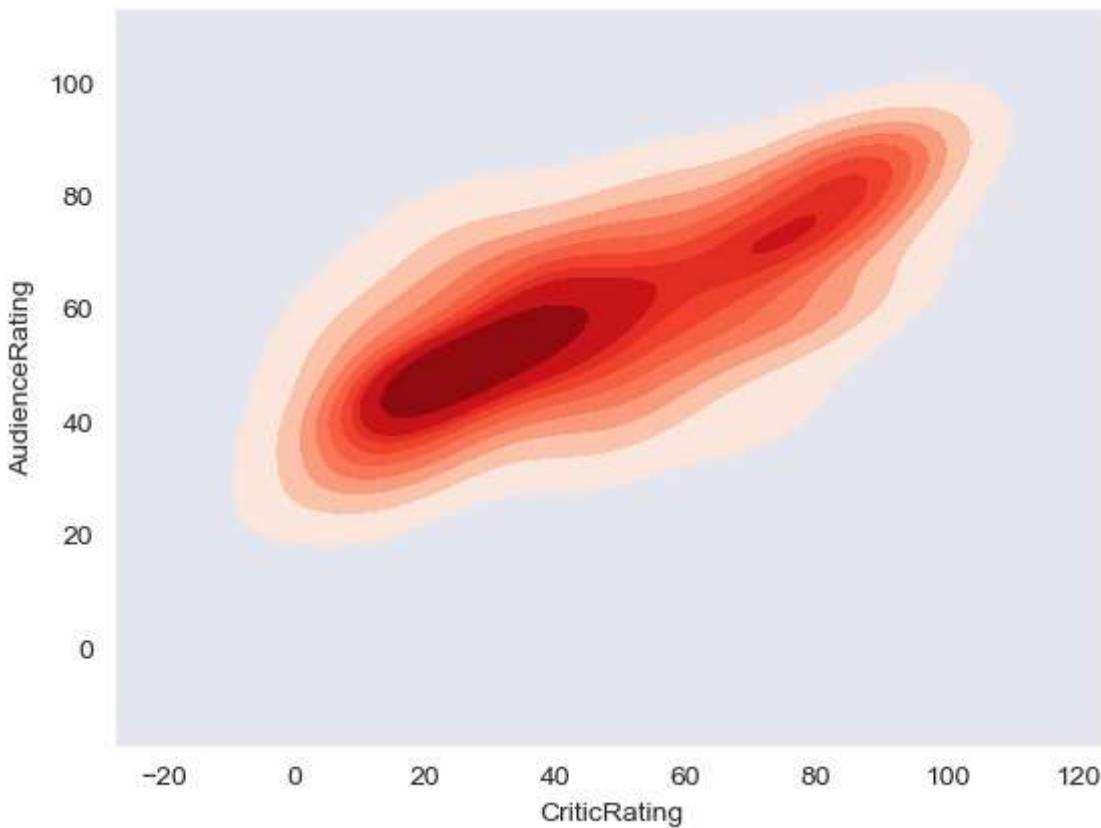
```
In [68]: # Kernel Density Estimate plot ( KDE PLOT )
# how can i visualize audience rating & critics rating . using scatterplot
```

```
In [70]: k1 = sns.kdeplot(data=movies,x='CriticRating',y='AudienceRating')

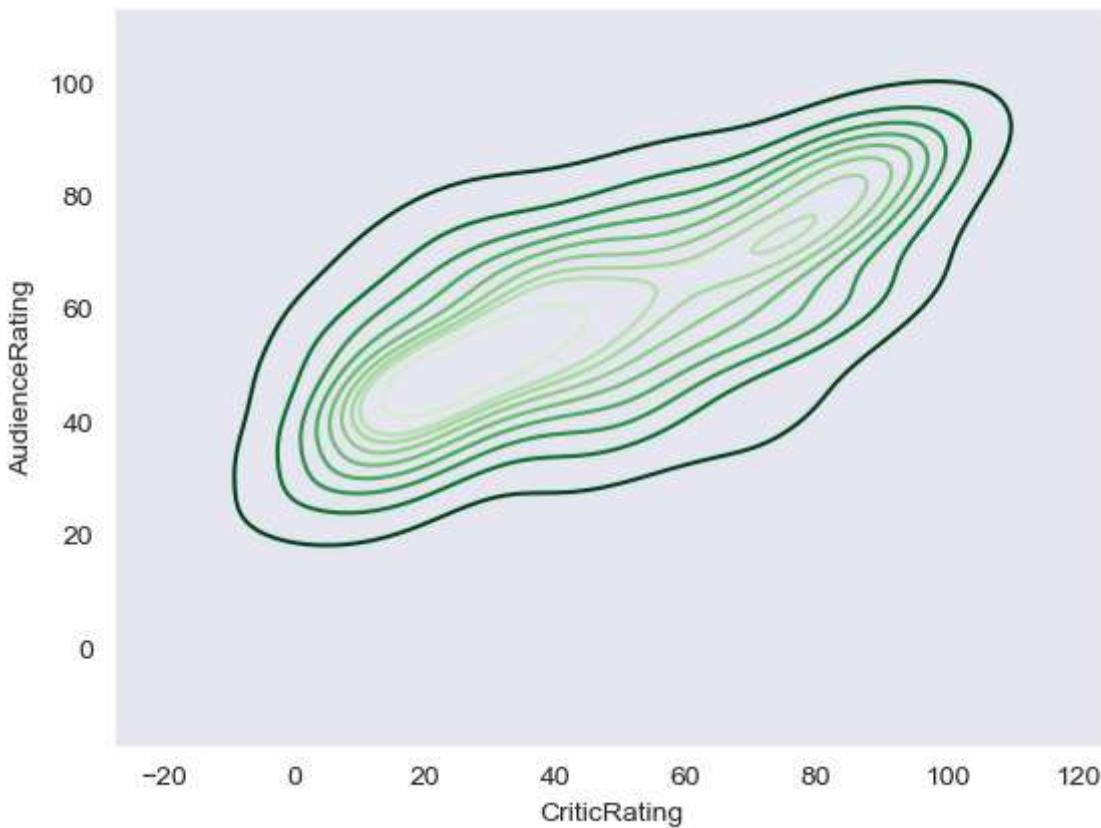
# where do u find more density and how density is distributed across from the the ch
# center point is kernel this is calld KDE & insteade of dots it visualize like thi
# we can able to clearly see the spread at the audience ratings
```



```
In [72]: k1 = sns.kdeplot(data=movies,x='CriticRating',y='AudienceRating',shade = True,shade
```

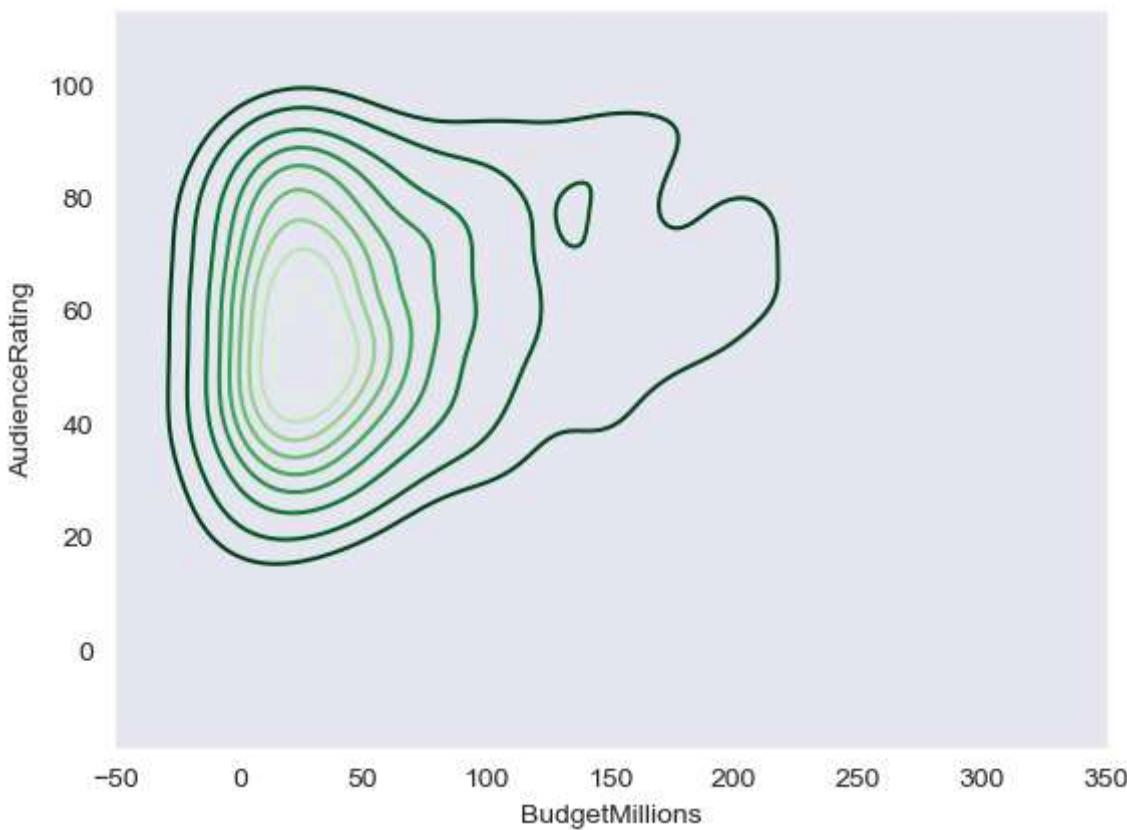


```
In [74]: k2 = sns.kdeplot(data=movies,x='CriticRating',y='AudienceRating',shade_lowest=False
```

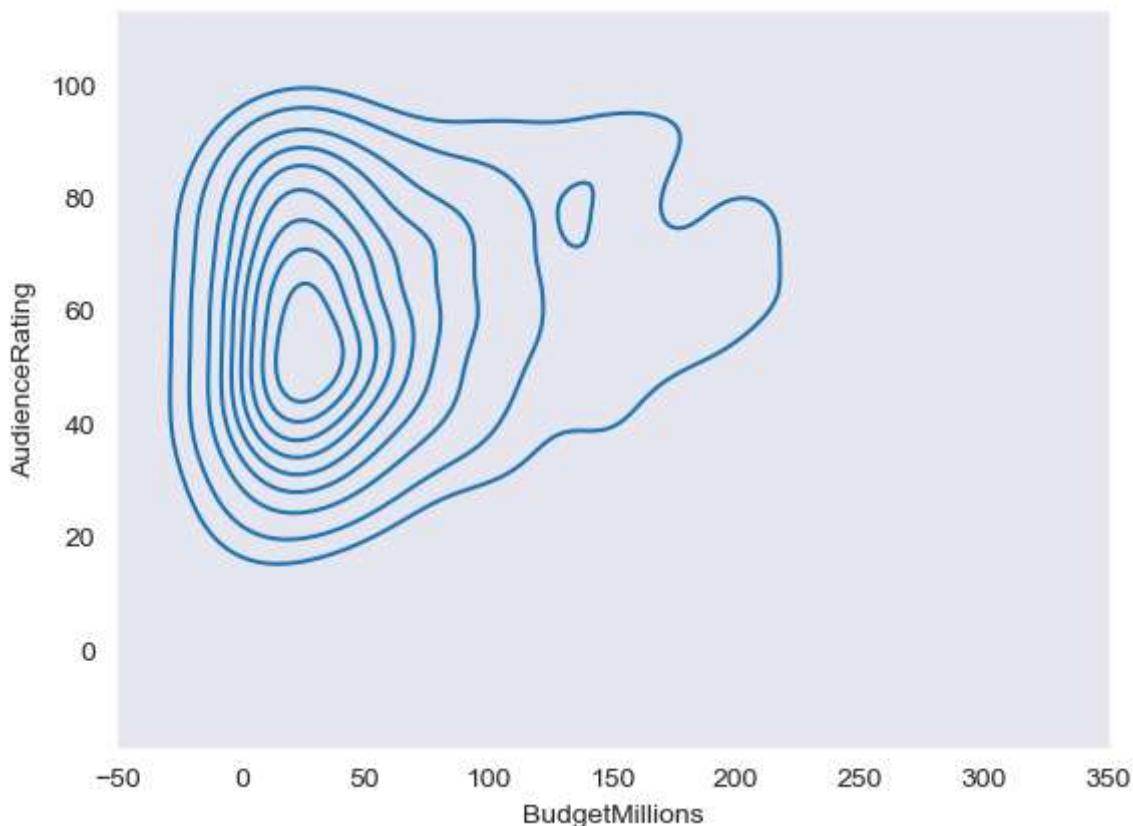


```
In [75]: sns.set_style('dark')
```

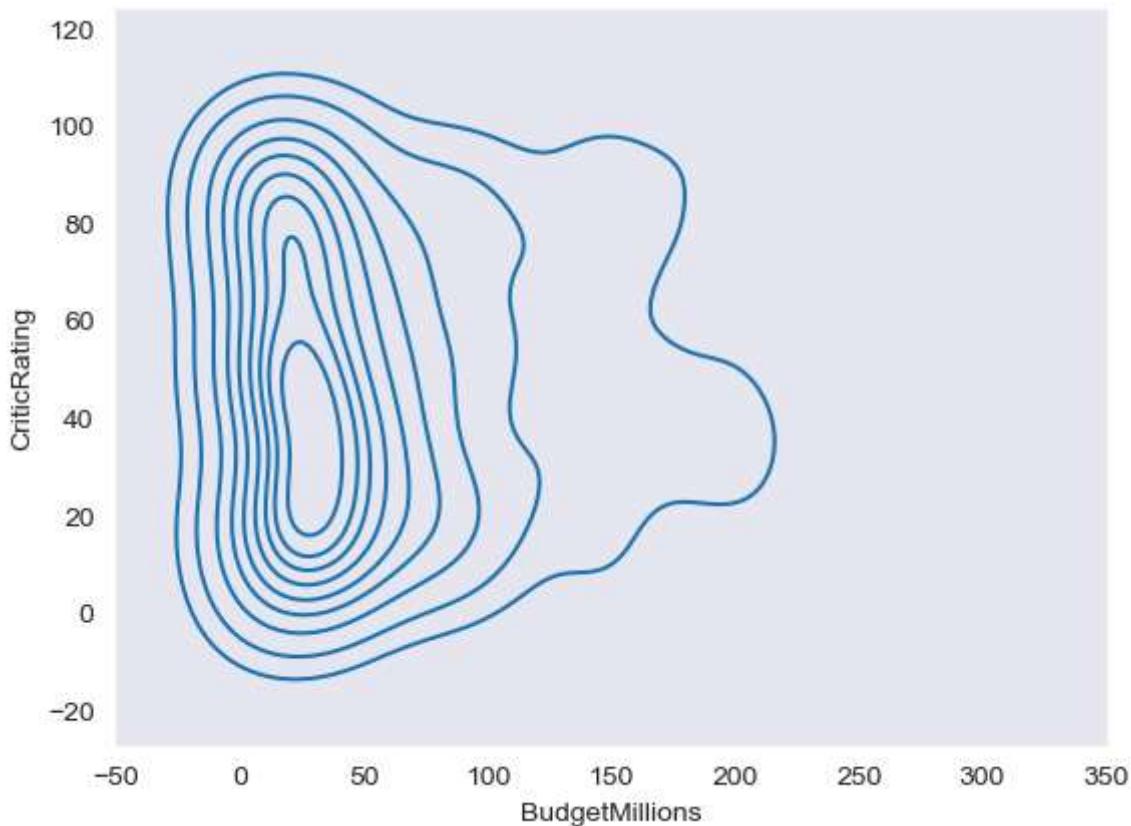
```
In [77]: k1 = sns.kdeplot(data=movies,x='BudgetMillions',y='AudienceRating',shade_lowest=False)
```



```
In [78]: sns.set_style('dark')
k1 = sns.kdeplot(data=movies,x='BudgetMillions',y='AudienceRating')
```



```
In [80]: sns.set_style('dark')
k2 = sns.kdeplot(data=movies,x='BudgetMillions',y='CriticRating')
```



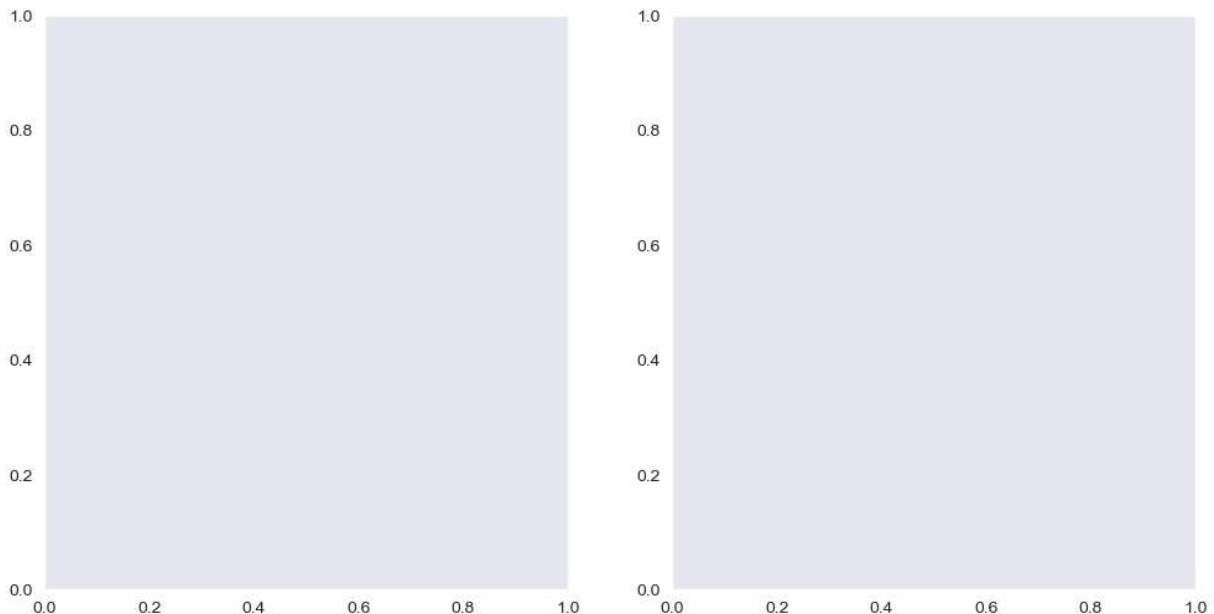
```
In [81]: movies.head()
```

```
Out[81]:
```

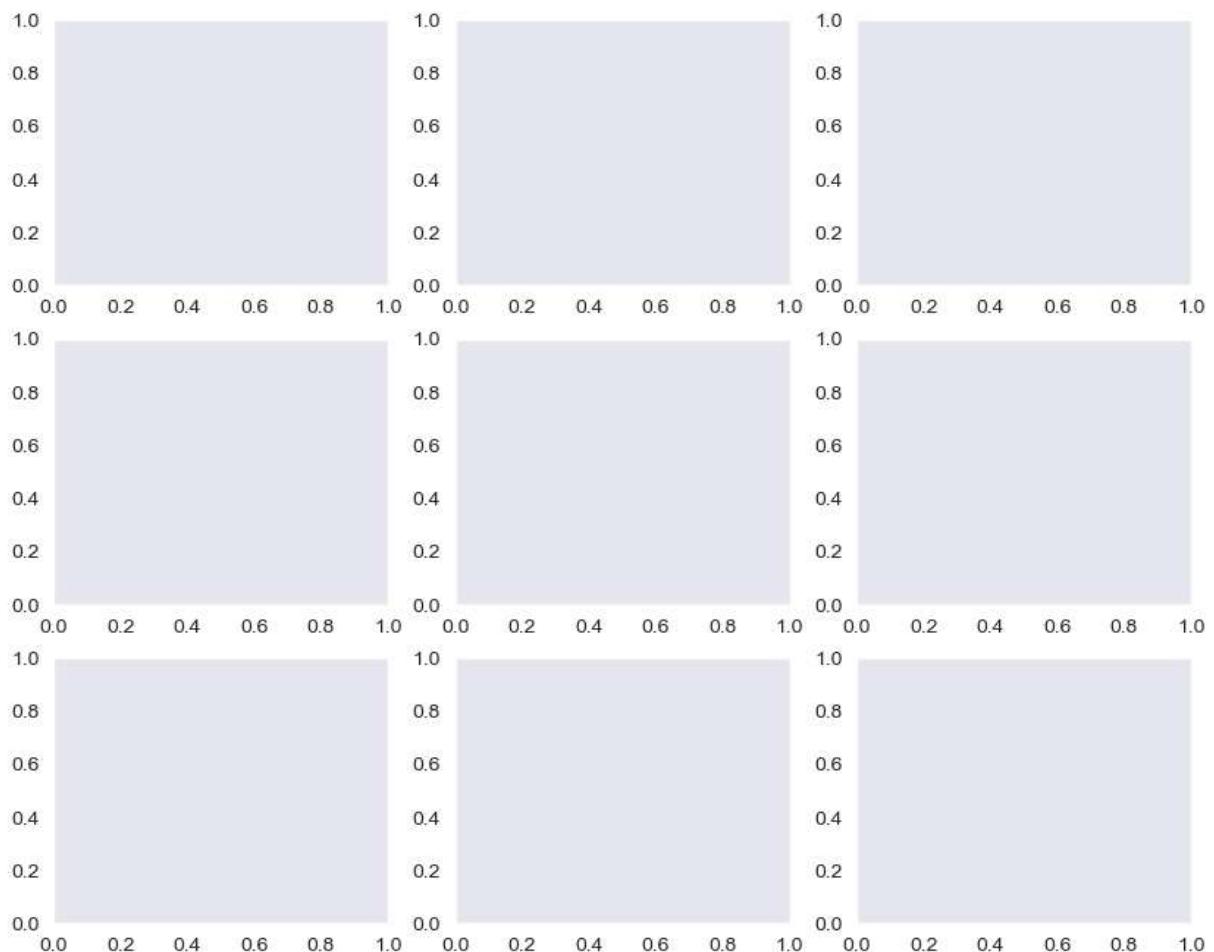
	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
0	(500) Days of Summer	Comedy	87	81	8	2009
1	10,000 B.C.	Adventure	9	44	105	2008
2	12 Rounds	Action	30	52	20	2009
3	127 Hours	Adventure	93	84	18	2010
4	17 Again	Comedy	55	70	20	2009

```
In [83]: #subplots
```

```
fig, ax = plt.subplots(1,2, figsize =(12,6))  
#f, ax = plt.subplots(3,3, figsize =(12,6))
```



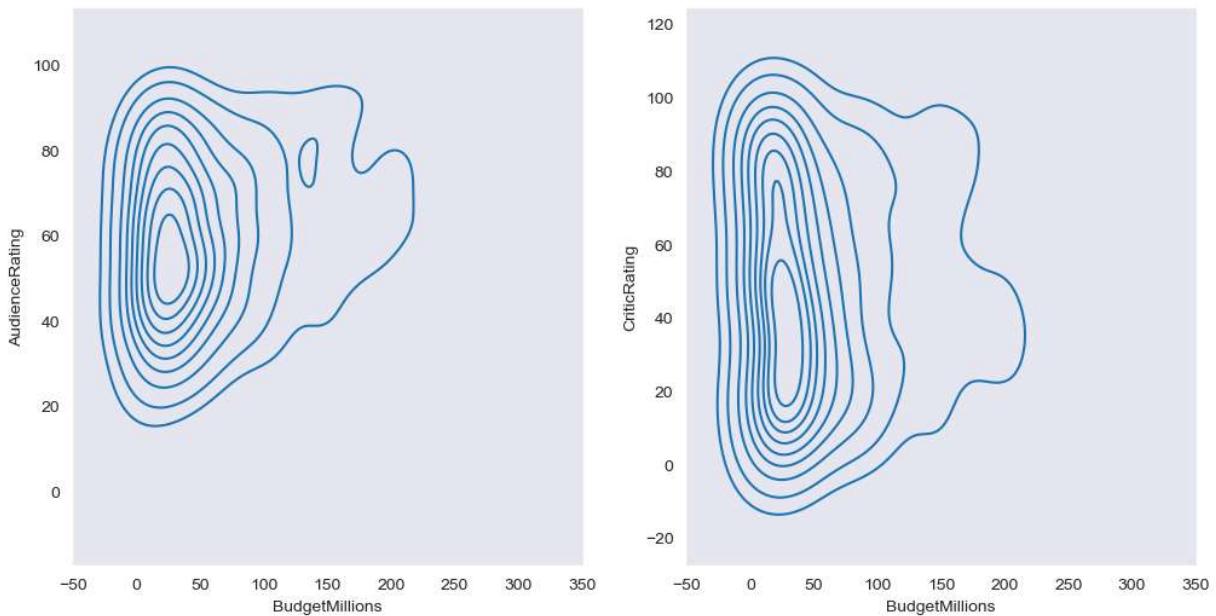
```
In [84]: #ax=plt.subplots(1,2,figsize=(3,3))
ax=plt.subplots(3,3,figsize=(10,8))
```



```
In [93]: # Create subplots
f, axes = plt.subplots(1, 2, figsize=(12, 6))

# Fix kdeplot usage by using keyword arguments for x and y
```

```
k1=sns.kdeplot(data=movies, x="BudgetMillions", y="AudienceRating", ax=axes[0])
k2=sns.kdeplot(data=movies, x="BudgetMillions", y="CriticRating", ax=axes[1])
```

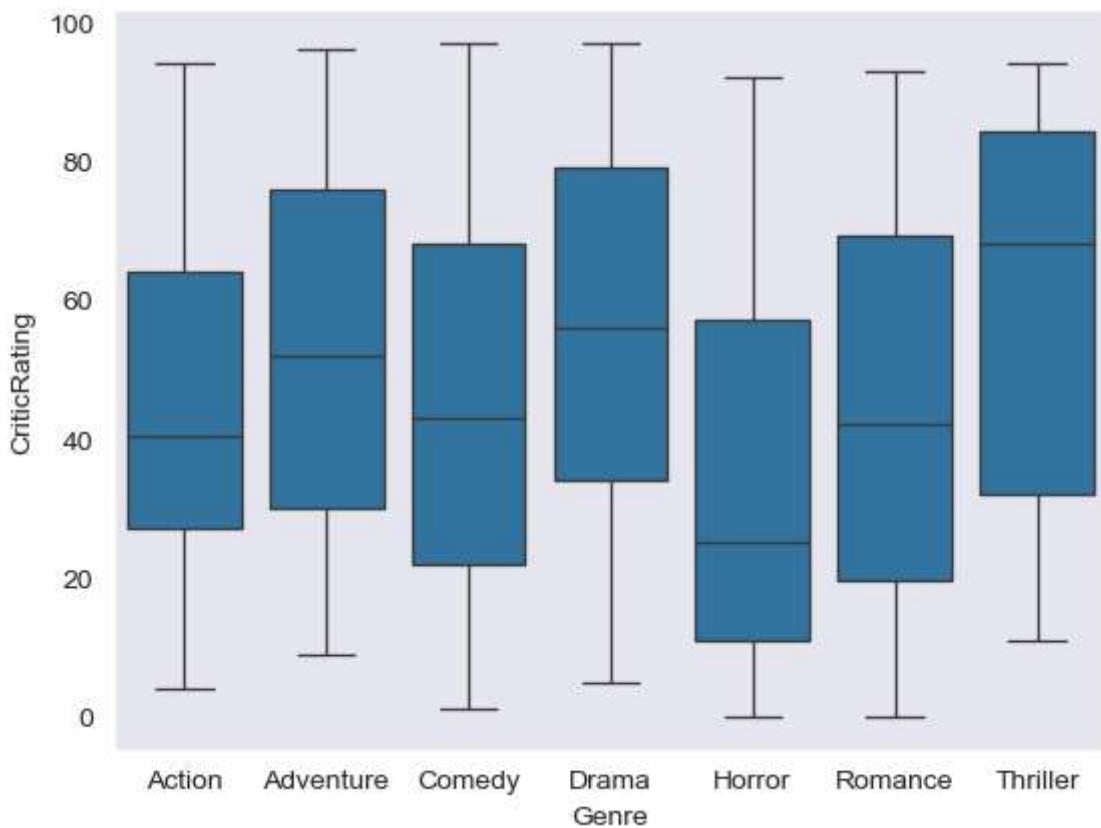


```
In [94]: axes
```

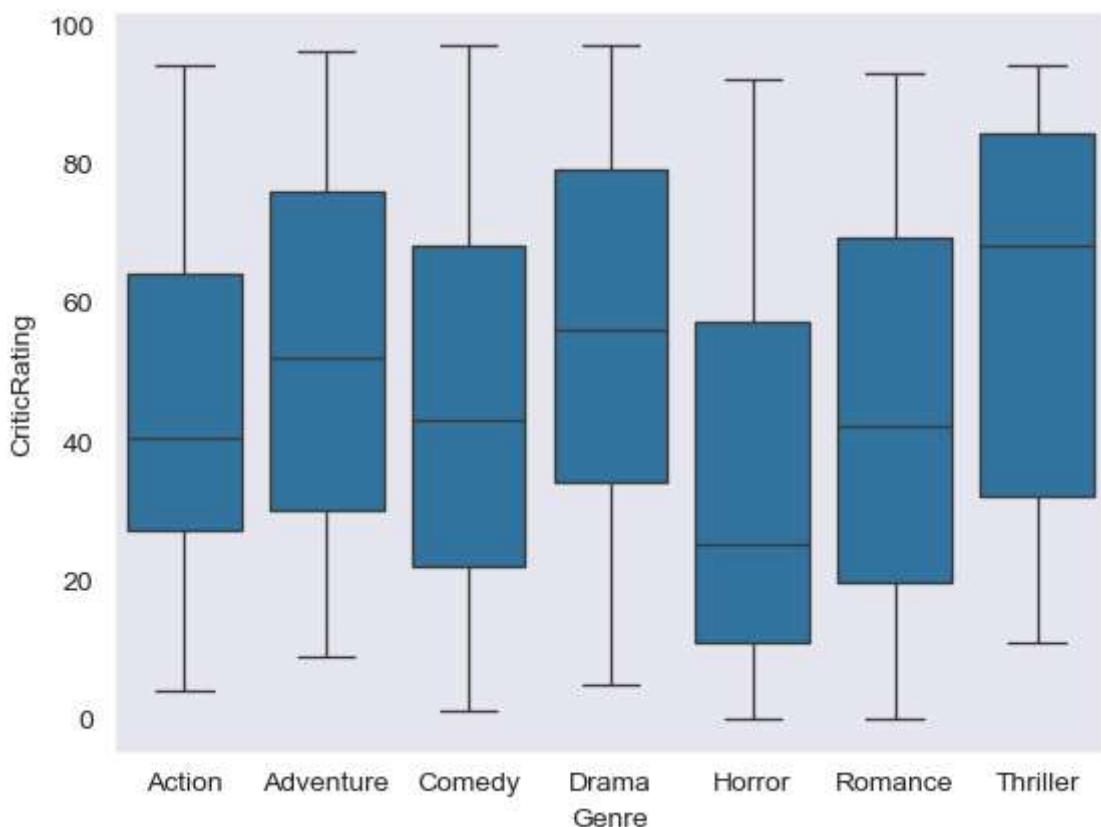
```
Out[94]: array([<Axes: xlabel='BudgetMillions', ylabel='AudienceRating'>,
   <Axes: xlabel='BudgetMillions', ylabel='CriticRating'>],
  dtype=object)
```

```
In [95]: #Box plots -
```

```
w = sns.boxplot(data=movies, x='Genre', y = 'CriticRating')
```

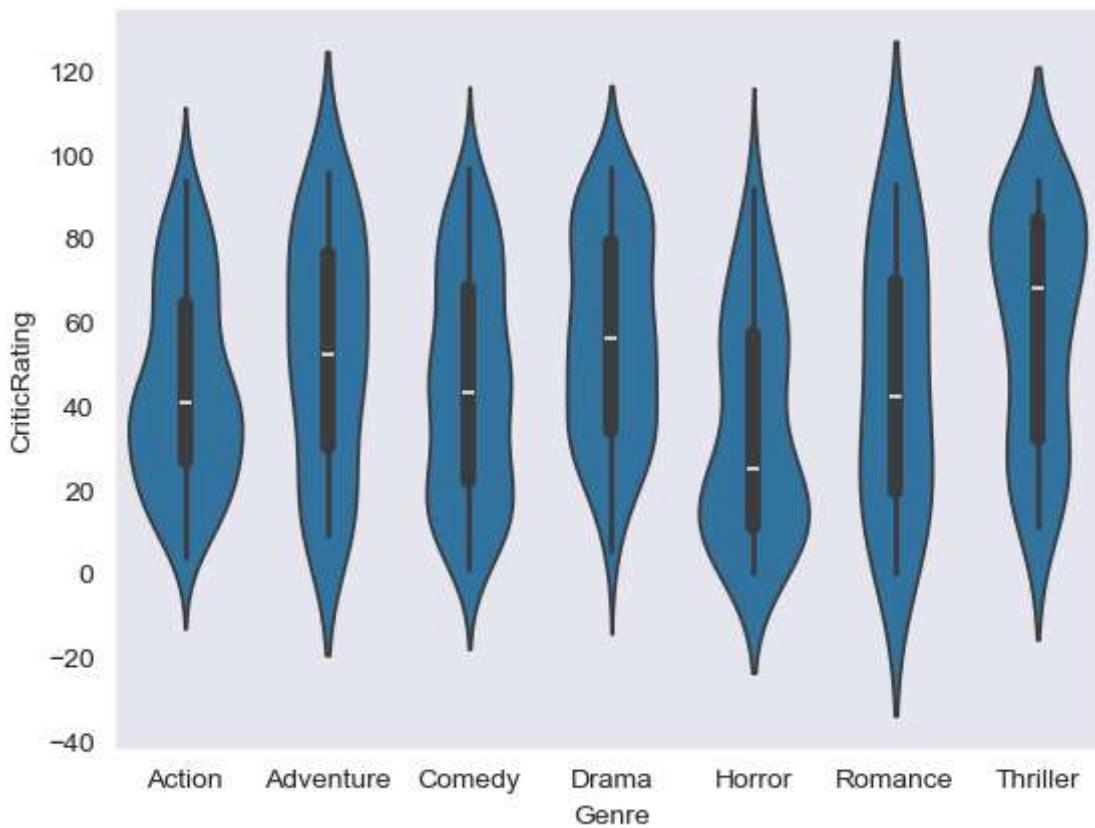


```
In [96]: sns.boxplot(data=movies, x='Genre', y='CriticRating')
```

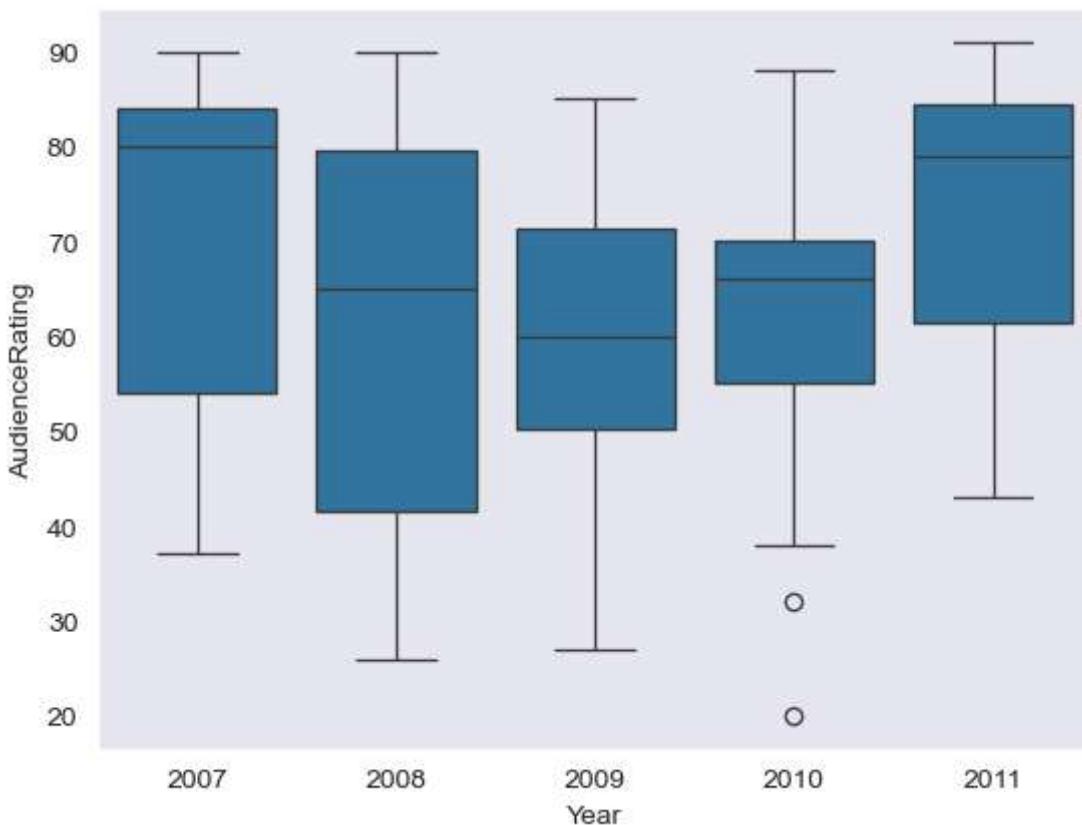


```
In [97]: #violin plot
```

```
z = sns.violinplot(data=movies, x='Genre', y = 'CriticRating')
```



```
In [98]: w1=sns.boxplot(data=movies[movies.Genre=='Drama'], x='Year',y='AudienceRating')
```



```
In [99]: movies.head()
```

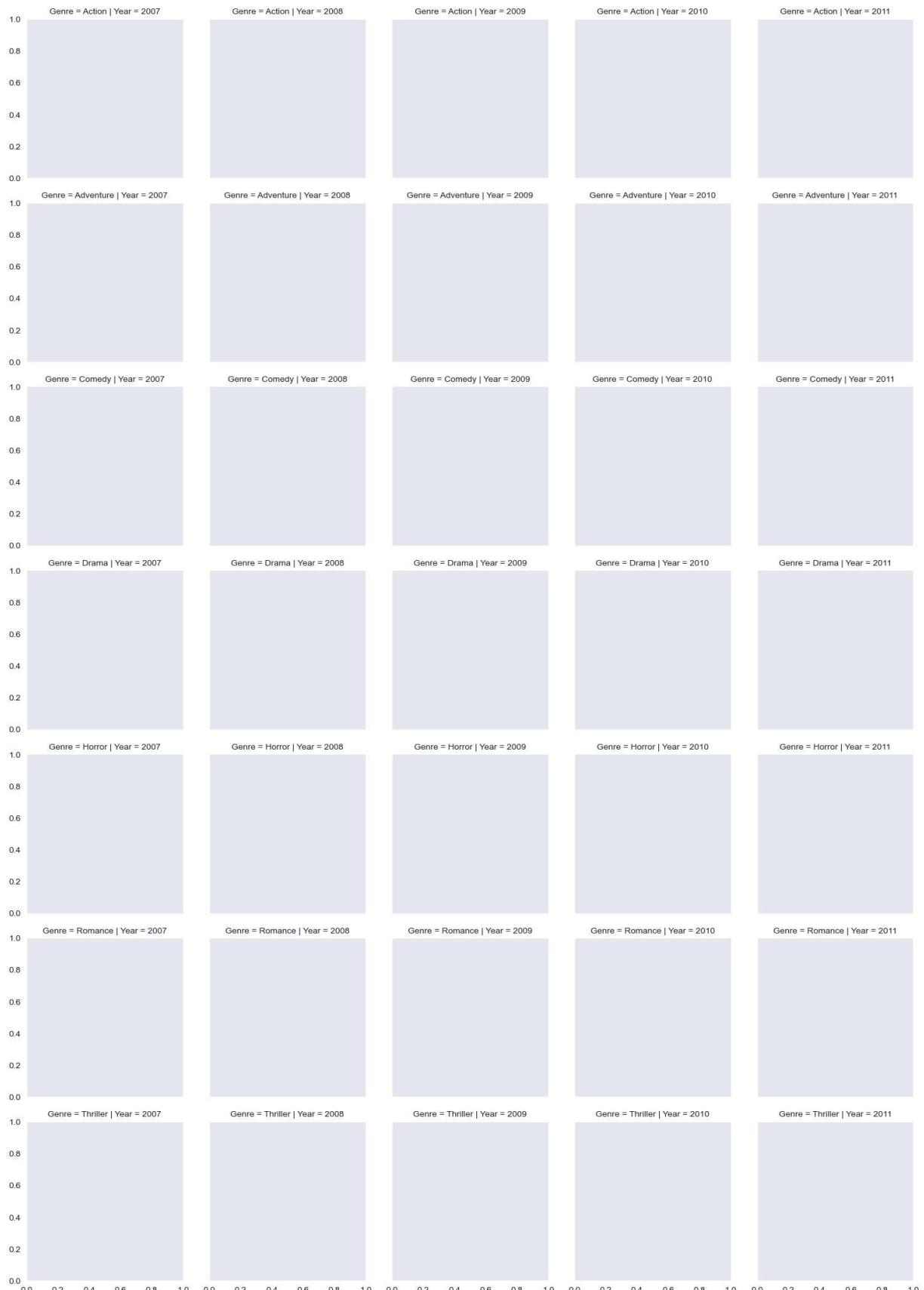
Out[99]:

	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
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 [102...]: # Createing a Facet grid
```

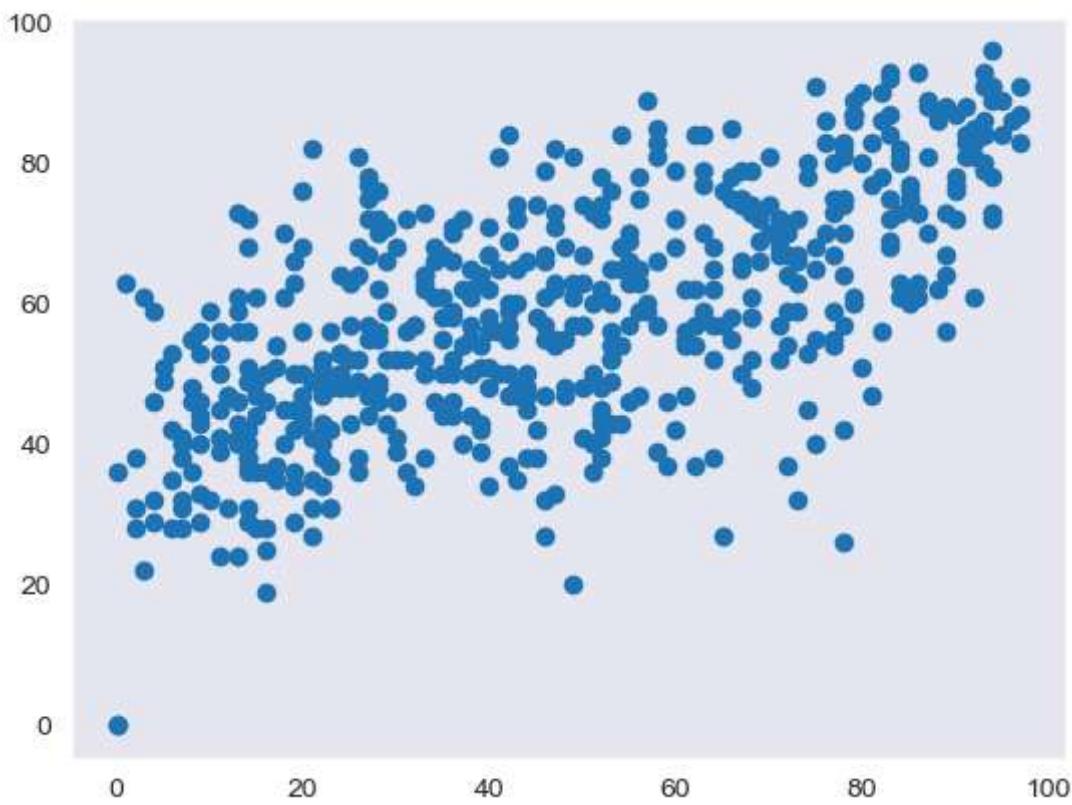
```
In [103...]: g =sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre') #kind of subp
```

Advanced_visualization_Movie_Ratings



```
In [105]: plt.scatter(movies.CriticRating, movies.AudienceRating)
```

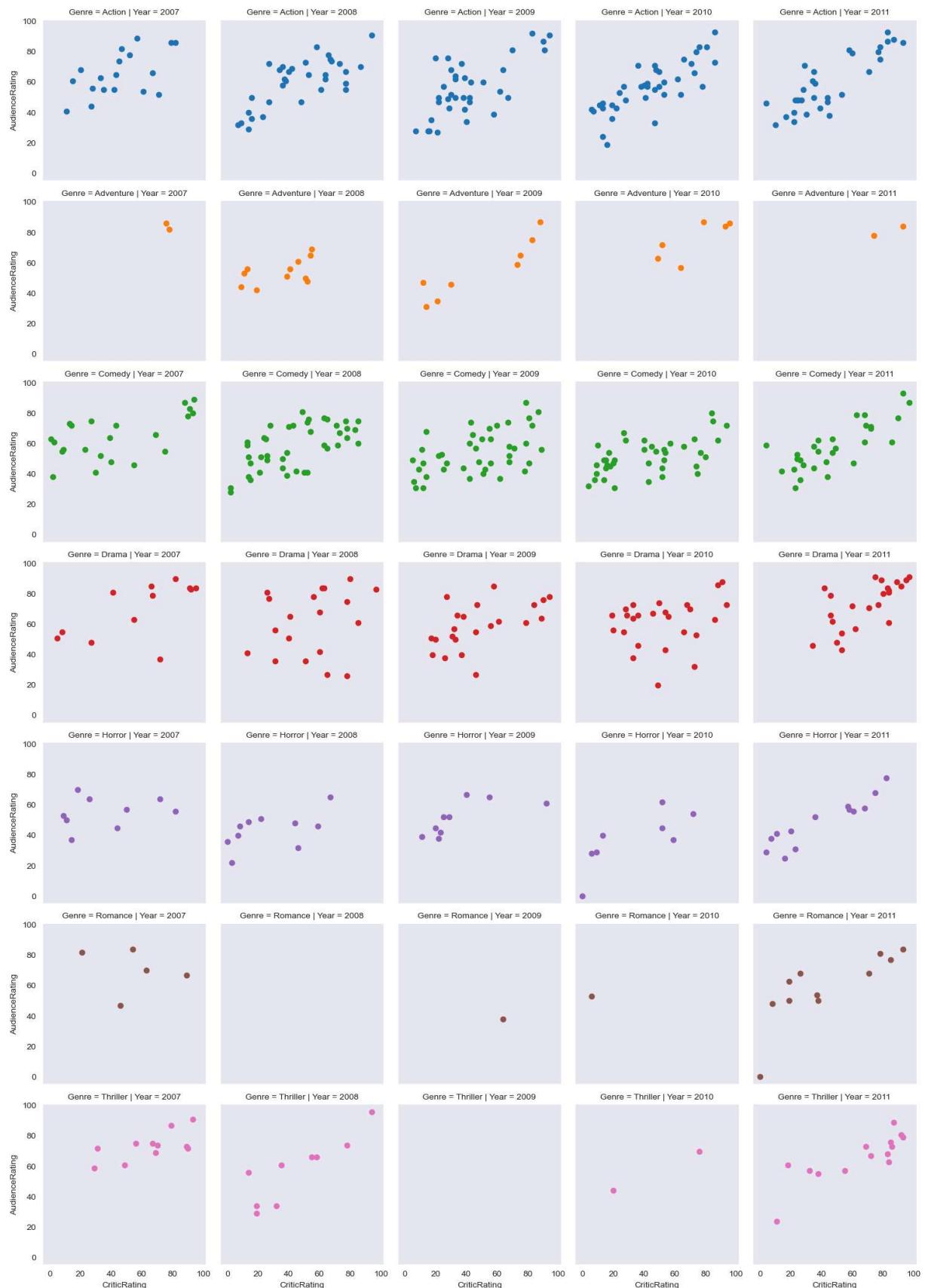
```
Out[105]: <matplotlib.collections.PathCollection at 0x20071499460>
```



In [106]:

```
g = sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre')
g = g.map(plt.scatter, 'CriticRating', 'AudienceRating' ) #scatterplots are mapped
```

Advanced_visualization_Movie_Ratings



```
In [107]: # you can populated any type of chat.
```

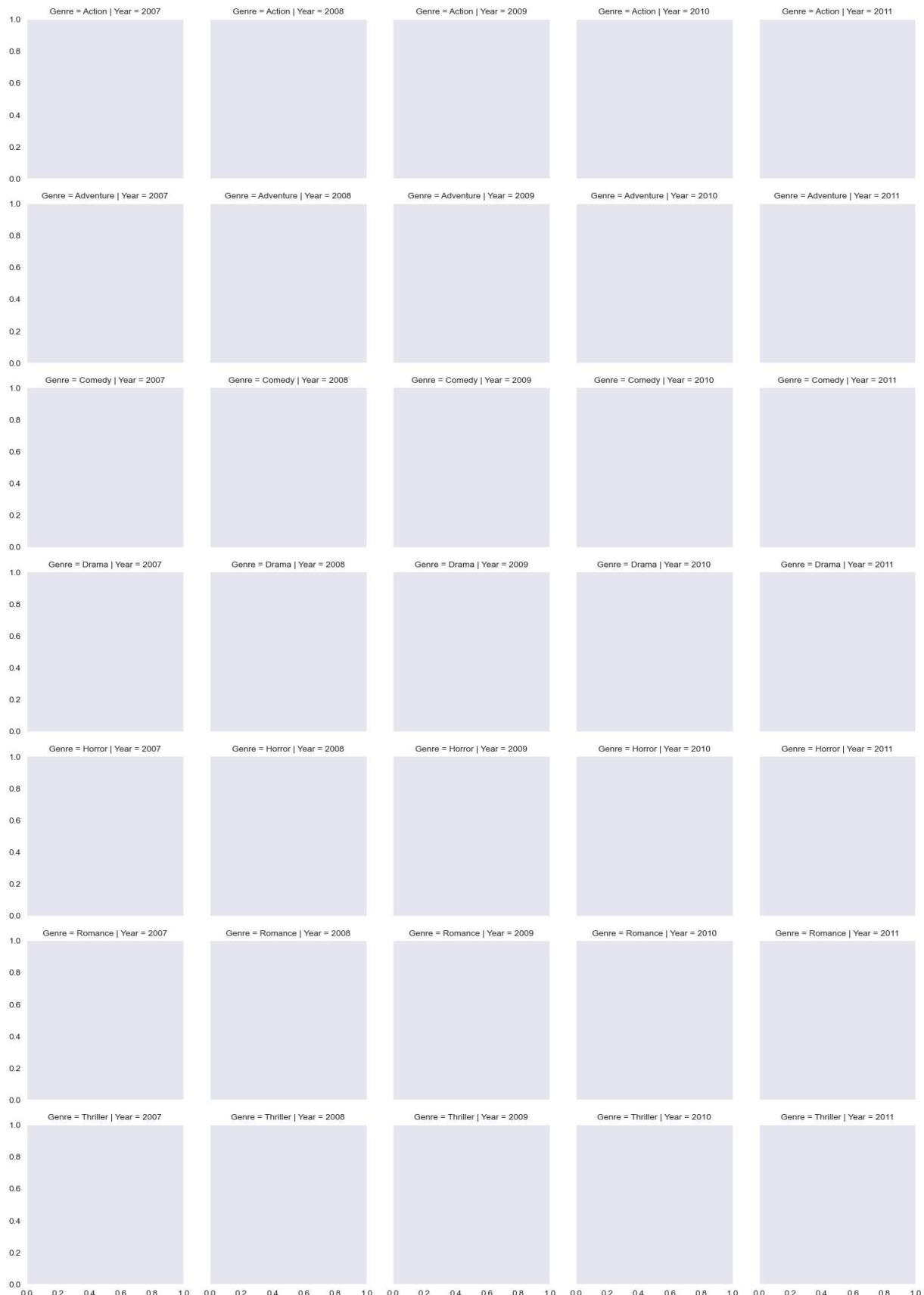
```
g = sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre')
g = g.map(plt.hist, 'BudgetMillions') #scatterplots are mapped in facetgrid
```

Advanced_visualization_Movie_Ratings



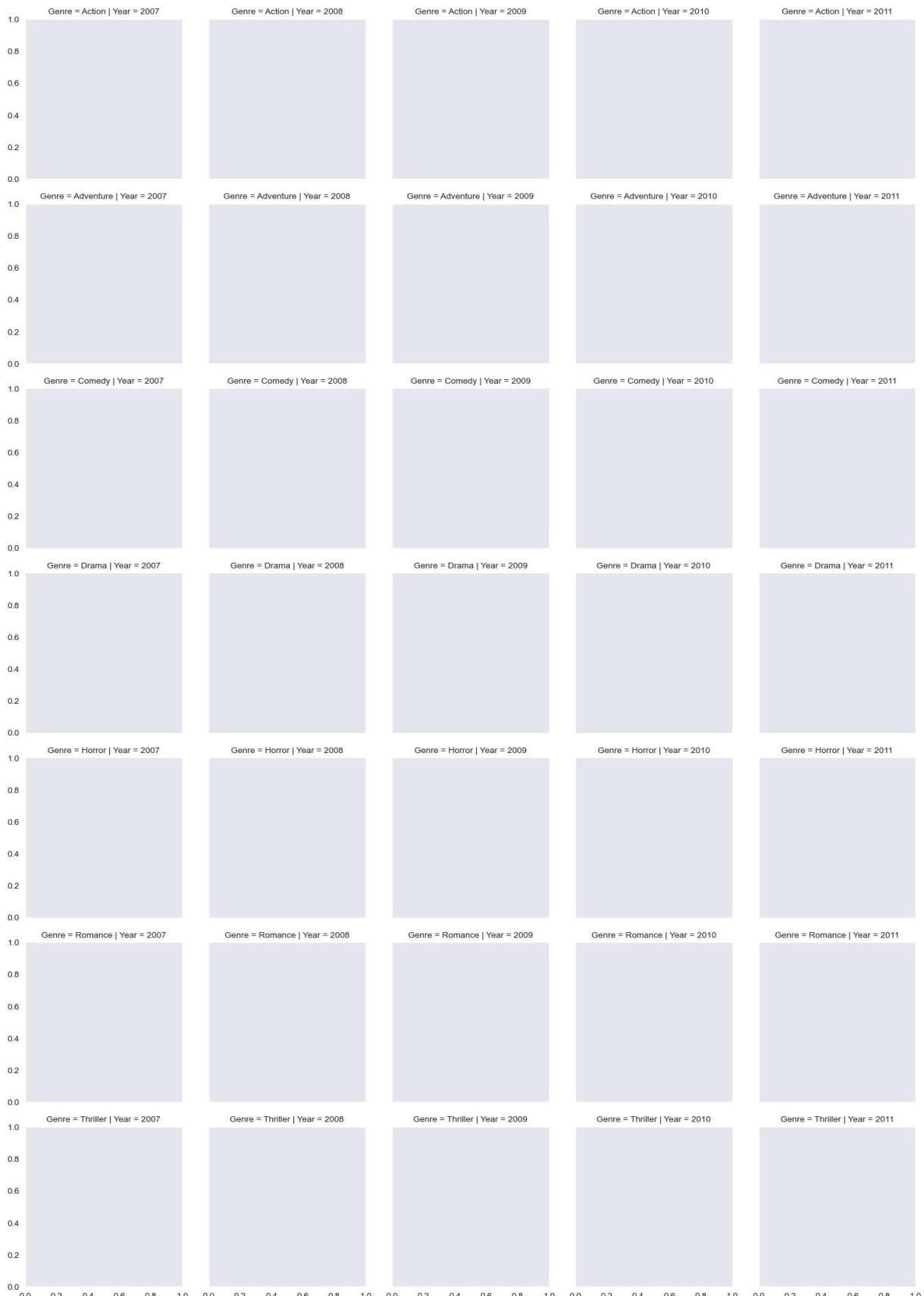
```
In [108...]: g = sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre')
#step 1
```

Advanced_visualization_Movie_Ratings



```
In [109]: g = sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre')
kws = dict(s=50, linewidth=0.5, edgecolor='black')
#step 1 & 2
```

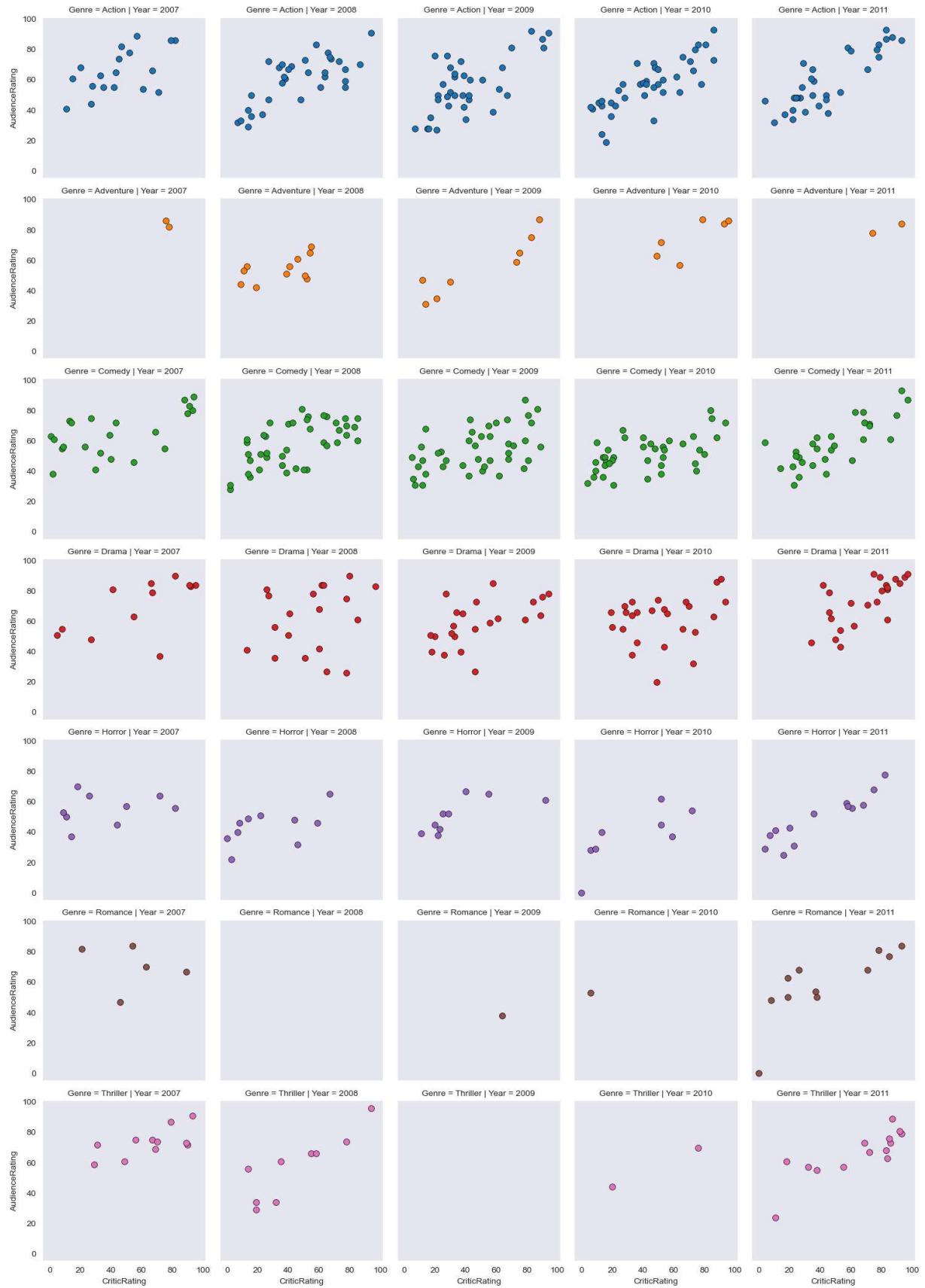
Advanced_visualization_Movie_Ratings



In [111...]

```
#  
g =sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre')  
kws = dict(s=50, linewidth=0.5,edgecolor='black')
```

```
g = g.map(plt.scatter, 'CriticRating', 'AudienceRating', **kws ) #scatterplots are made
#step 1&2&3
```



In [115]:

```
# python is not vectorize programming Language
# Building dashboards (dashboard - combination of chats)
```

```
import seaborn as sns
import matplotlib.pyplot as plt

sns.set_style("darkgrid") # Set dark grid style

# Create subplots
f, axes = plt.subplots(2, 2, figsize=(15, 15))

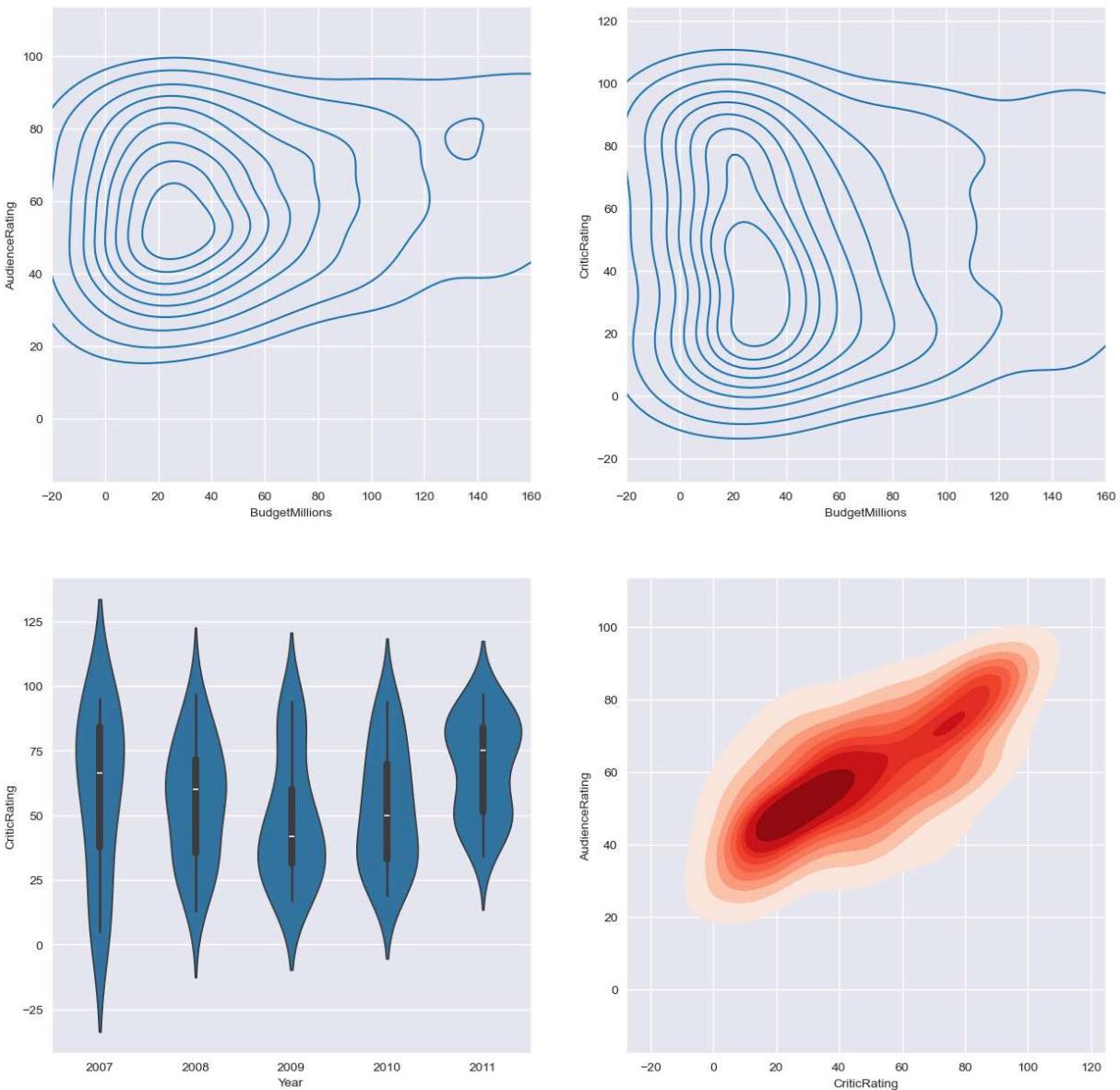
# KDE plots
sns.kdeplot(data=movies, x="BudgetMillions", y="AudienceRating", ax=axes[0,0])
sns.kdeplot(data=movies, x="BudgetMillions", y="CriticRating", ax=axes[0,1])

# Set x-axis limits
axes[0,0].set(xlim=(-20,160))
axes[0,1].set(xlim=(-20,160))

# Violin plot for Drama genre
sns.violinplot(data=movies[movies.Genre=="Drama"], x="Year", y="CriticRating", ax=axes[1,0])

# KDE plot with shading
sns.kdeplot(data=movies, x="CriticRating", y="AudienceRating", shade=True, cmap="Reds")

plt.show()
```



In [114...]

```

import seaborn as sns
import matplotlib.pyplot as plt

sns.set_style("dark", {"axes.facecolor": "black"}) # Set dark background

# Create subplots
f, axes = plt.subplots(2, 2, figsize=(15, 15))

# KDE plots with corrected parameters
sns.kdeplot(data=movies, x="BudgetMillions", y="AudienceRating", shade=True, cmap="Blues")
sns.kdeplot(data=movies, x="BudgetMillions", y="AudienceRating", cmap="cool", ax=axes[0][1])

sns.kdeplot(data=movies, x="BudgetMillions", y="CriticRating", shade=True, cmap="inferno")
sns.kdeplot(data=movies, x="BudgetMillions", y="CriticRating", cmap="cool", ax=axes[1][0])

# Violin plot for Drama genre
sns.violinplot(data=movies[movies.Genre=="Drama"], x="Year", y="CriticRating", ax=axes[1][1])

# KDE plot with shading
sns.kdeplot(data=movies, x="CriticRating", y="AudienceRating", shade=True, cmap="Blues")

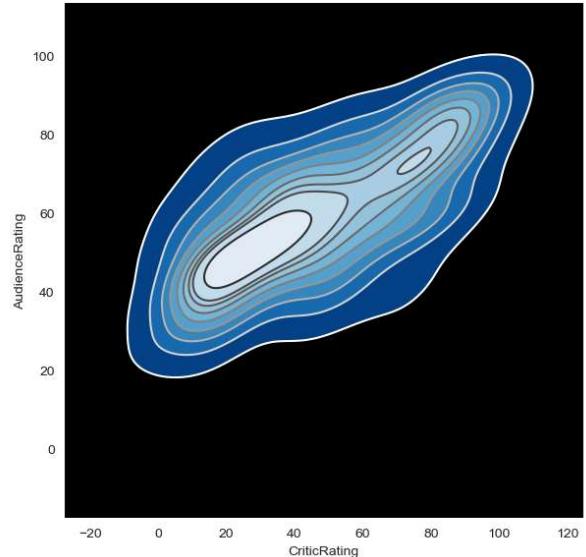
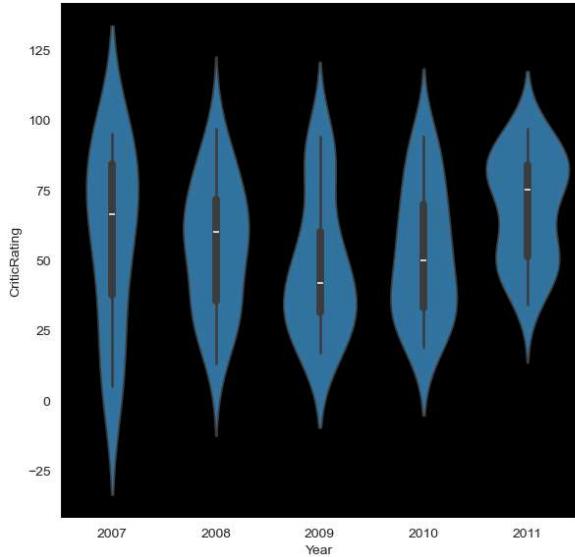
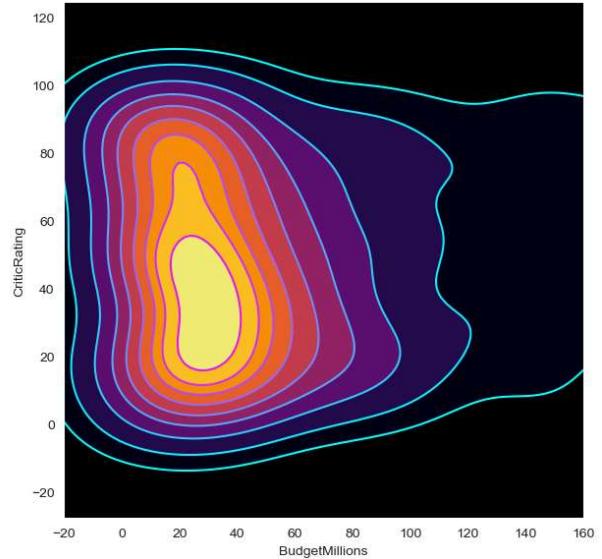
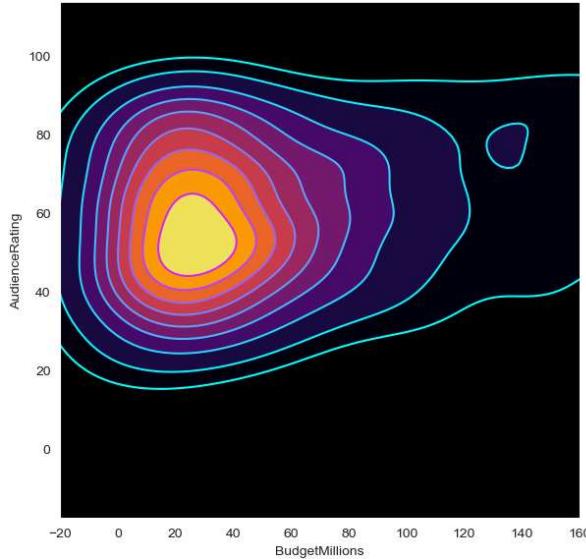
```

```

sns.kdeplot(data=movies, x="CriticRating", y="AudienceRating", cmap="gist_gray_r",
             # Set x-axis limits
             axes[0,0].set(xlim=(-20,160))
             axes[0,1].set(xlim=(-20,160))

plt.show()

```



Final discussion what we learn so far - 1> category datatype in python

2> jointplots

3> histogram

4> stacked histograms

5> Kde plot

6> subplot

- 7> violin plots
- 8> Facet grid
- 9> Building dashboards