

MOVIE RATING ANALYTICS (ADVANCED VISULIZATION) using SEABORN

```
In [13]: import pandas as pd
import os
```

```
In [14]: os.getcwd() # if you want to change the working directory
```

```
Out[14]: 'C:\\Users\\gadel'
```

```
In [15]: movies = pd.read_csv(r'C:\\Users\\gadel\\OneDrive\\Desktop\\Nareshit DataScience by Pro
```

```
In [16]: movies
```

```
Out[16]:
```

	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 [17]: len(movies)
```

```
Out[17]: 559
```

```
In [18]: movies.head()
```

Out[18]:

	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 [19]: `movies.tail()`

Out[19]:

	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 [20]: `movies.columns`

Out[20]: Index(['Film', 'Genre', 'Rotten Tomatoes Ratings %', 'Audience Ratings %', 'Budget (million \$)', 'Year of release'], dtype='object')

In [21]: `movies.columns = ['Film', 'Genre', 'CriticRating', 'AudienceRating', 'BudgetMillion`In [22]: `movies.head() # Removed spaces & % removed noise characters`

Out[22]:

	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 [23]: `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 [24]: 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[24]:
```

	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 [25]: movies['Film']
#movies['Audience Ratings %']
```

```
Out[25]: 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: object
```

```
In [26]: movies['Film']
#movies['Audience Ratings %']
```

```
Out[26]: 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: object
```

```
In [27]: movies.Film = movies.Film.astype('category')
```

```
In [28]: movies.Film
```

```
Out[28]: 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 [29]: movies.head()
```

```
Out[29]:
```

	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 [30]: movies.info()

# now the same thing we will change genra to category & year to category
```

```
<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 [31]: movies.Genre = movies.Genre.astype('category')
movies.Year = movies.Year.astype('category')
```

```
In [32]: movies.Genre
```

```
Out[32]: 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 [33]: movies.Year # is it real no. year you can take average,min,max but out come have r
```

```
Out[33]: 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 [34]: 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 [35]: movies.Genre.cat.categories
```

```
Out[35]: Index(['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance',
               'Thriller'],
              dtype='object')
```

```
In [36]: movies.describe()
#now when you see the descript you will get only integer value mean, standard dev
```

```
Out[36]:
```

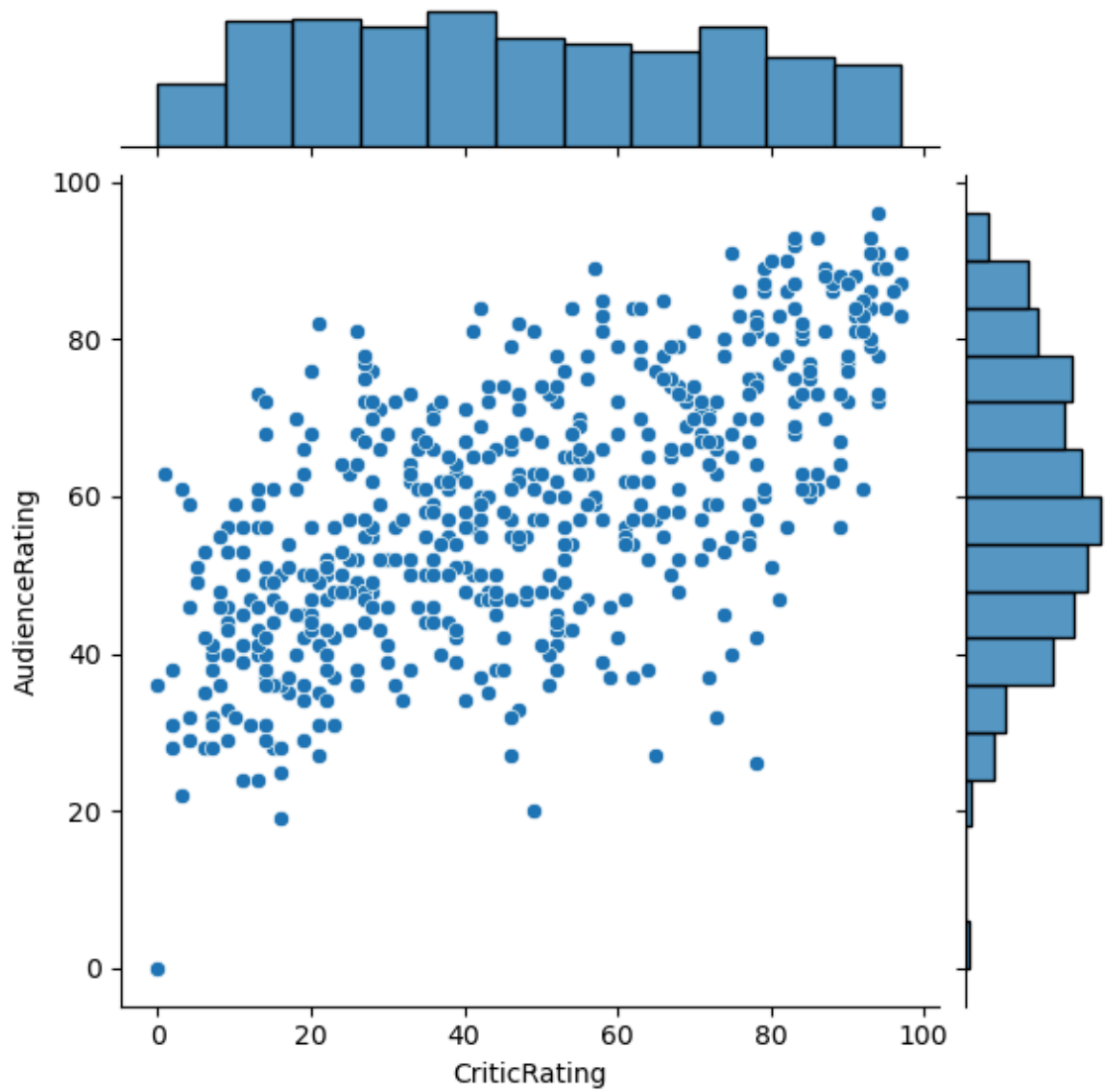
	CriticRating	AudienceRating	BudgetMillions
count	559.000000	559.000000	559.000000
mean	47.309481	58.744186	50.236136
std	26.413091	16.826887	48.731817
min	0.000000	0.000000	0.000000
25%	25.000000	47.000000	20.000000
50%	46.000000	58.000000	35.000000
75%	70.000000	72.000000	65.000000
max	97.000000	96.000000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

```
In [37]: # How to working with joint plots

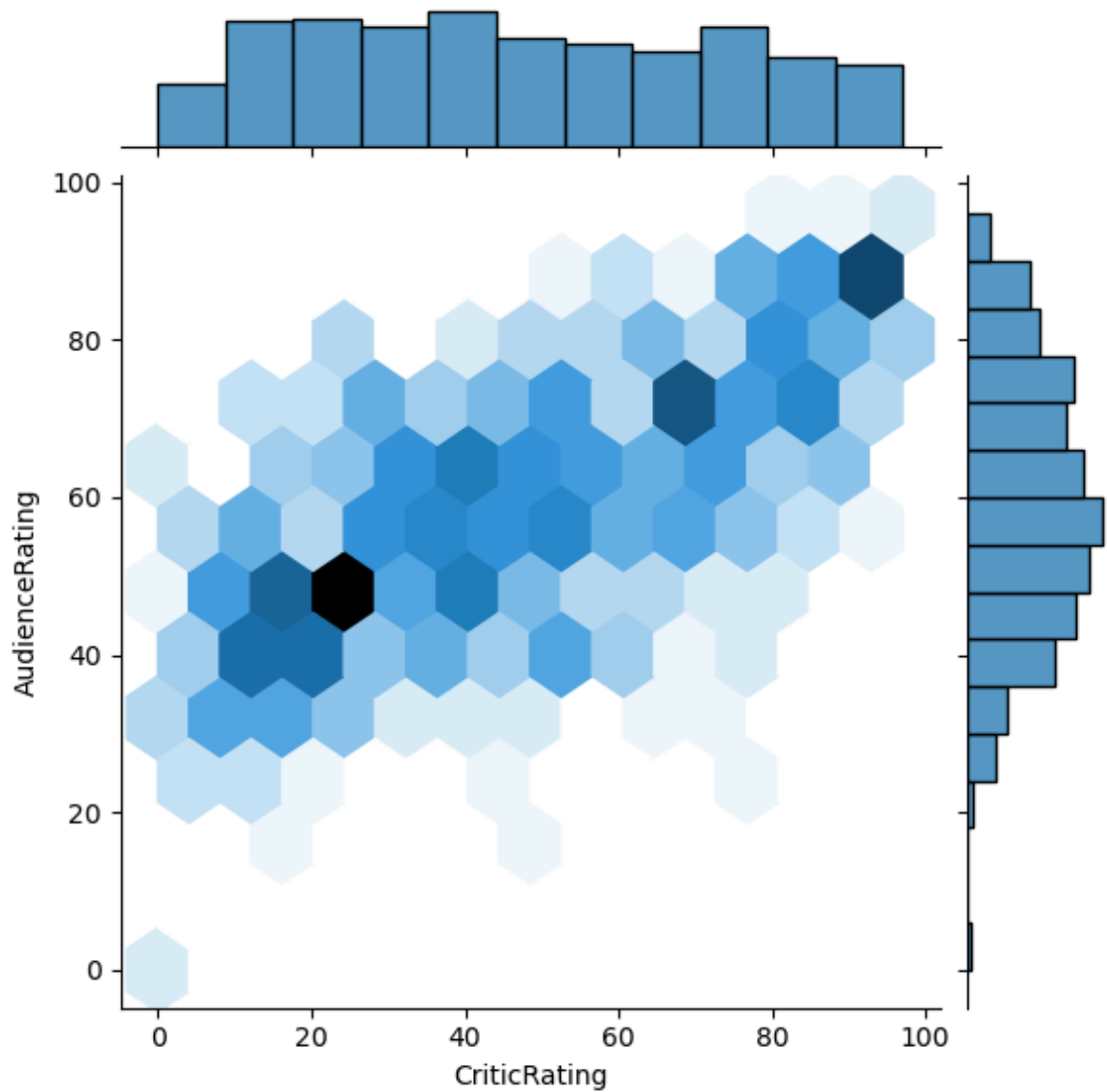
from matplotlib import pyplot as plt
import seaborn as sns
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
```

```
In [ ]: * basically joint plot is a scatter plot & it find the relation b/w audiene & crit
        * also if you look up you can find the uniform distribution (critics)and normal di
```

```
In [39]: j = sns.jointplot( data = movies, x = 'CriticRating', y = 'AudienceRating')
# Audience rating is more dominant then critics rating
# Based on this we find out as most people are most liklihood to watch audience re
# Let me explain the excel - if you filter audience rating & critic rating. critic
```



```
In [42]: j = sns.jointplot( data = movies, x = 'CriticRating', y = 'AudienceRating', kind='
#j = sns.jointplot( data = movies, x = 'CriticRating', y = 'AudienceRating', kind=
```

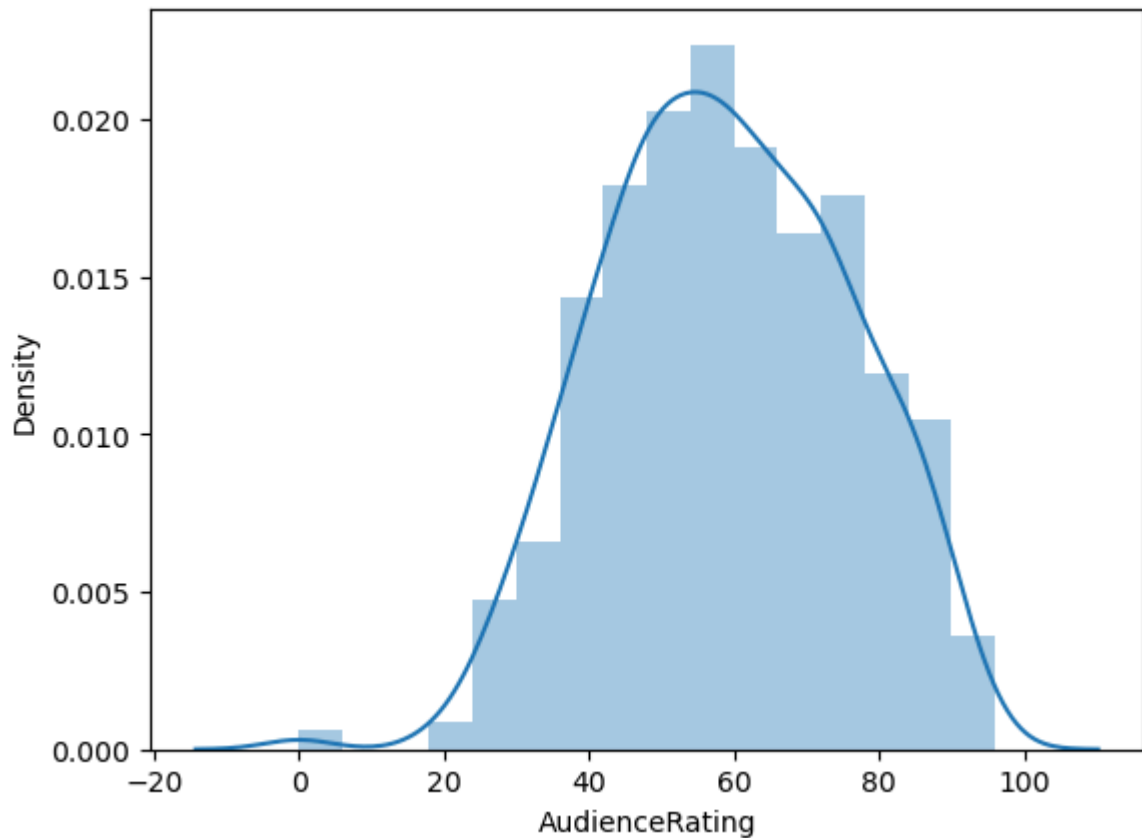


```
In [44]: #Histograms

# <<< chat1

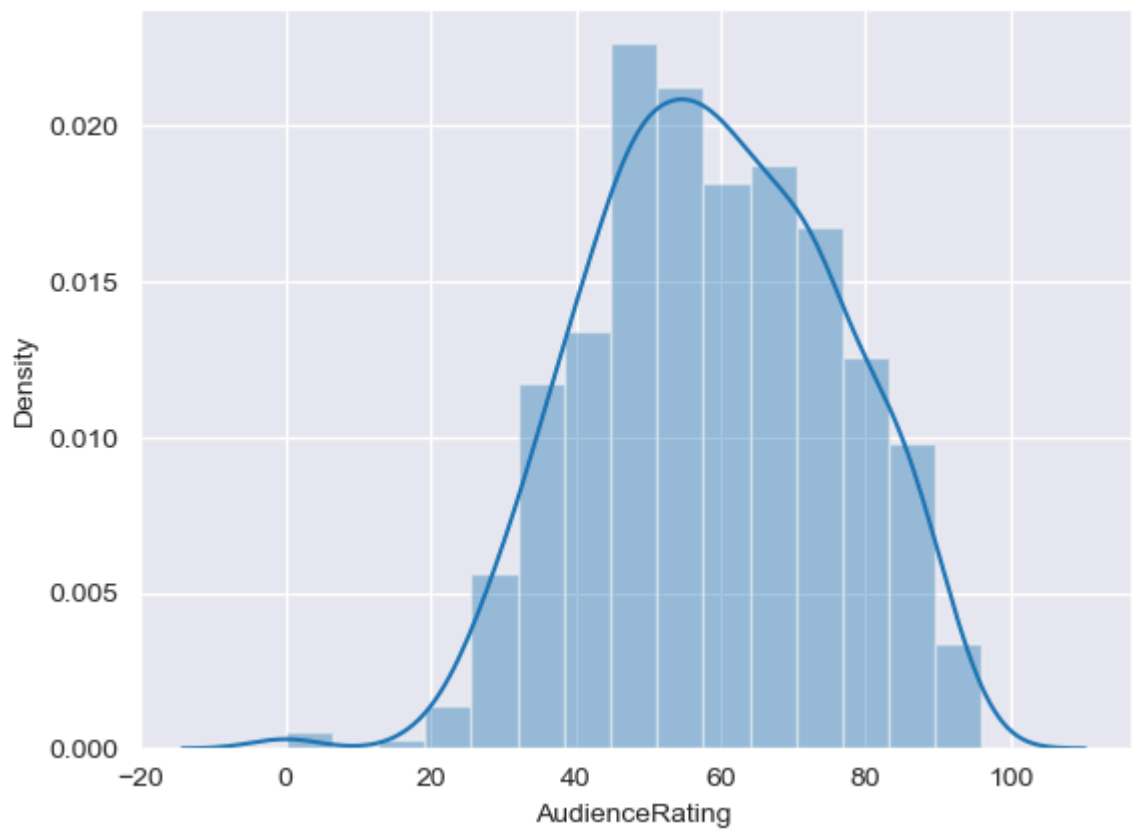
m1 = sns.distplot(movies.AudienceRating)

#y - axis generated by seaborn automatically that is the powerfull of seaborn galle
```

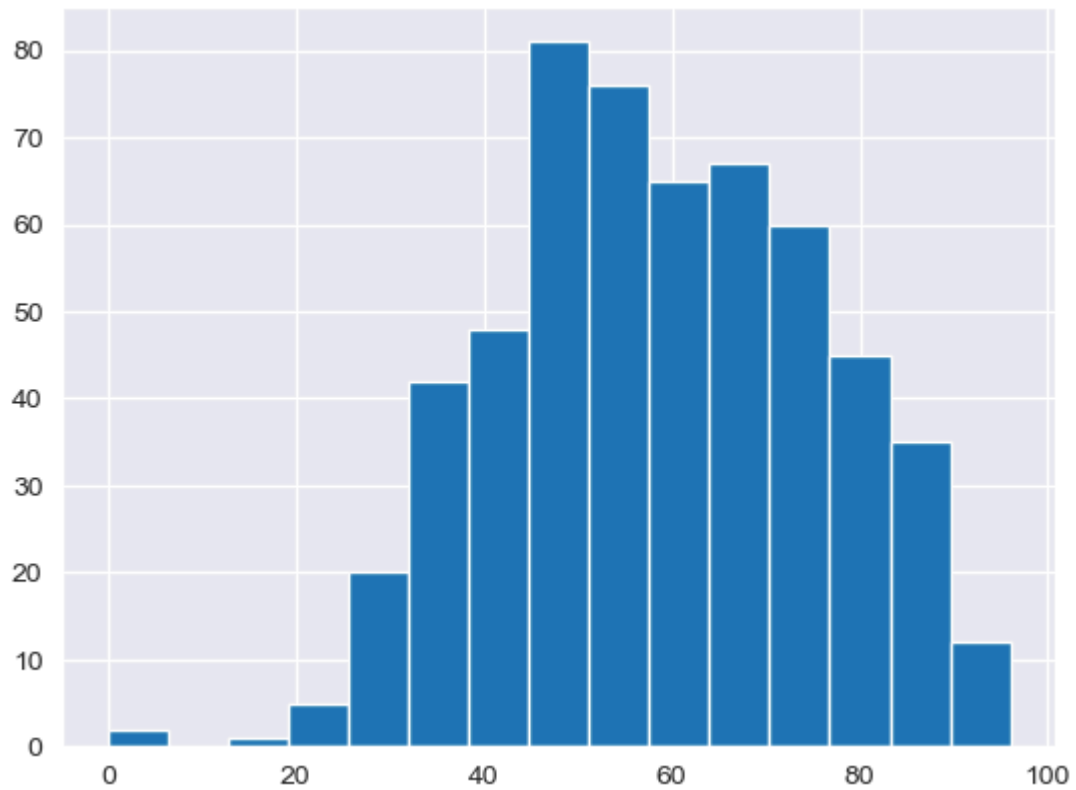



```
In [46]: sns.set_style('darkgrid')
```

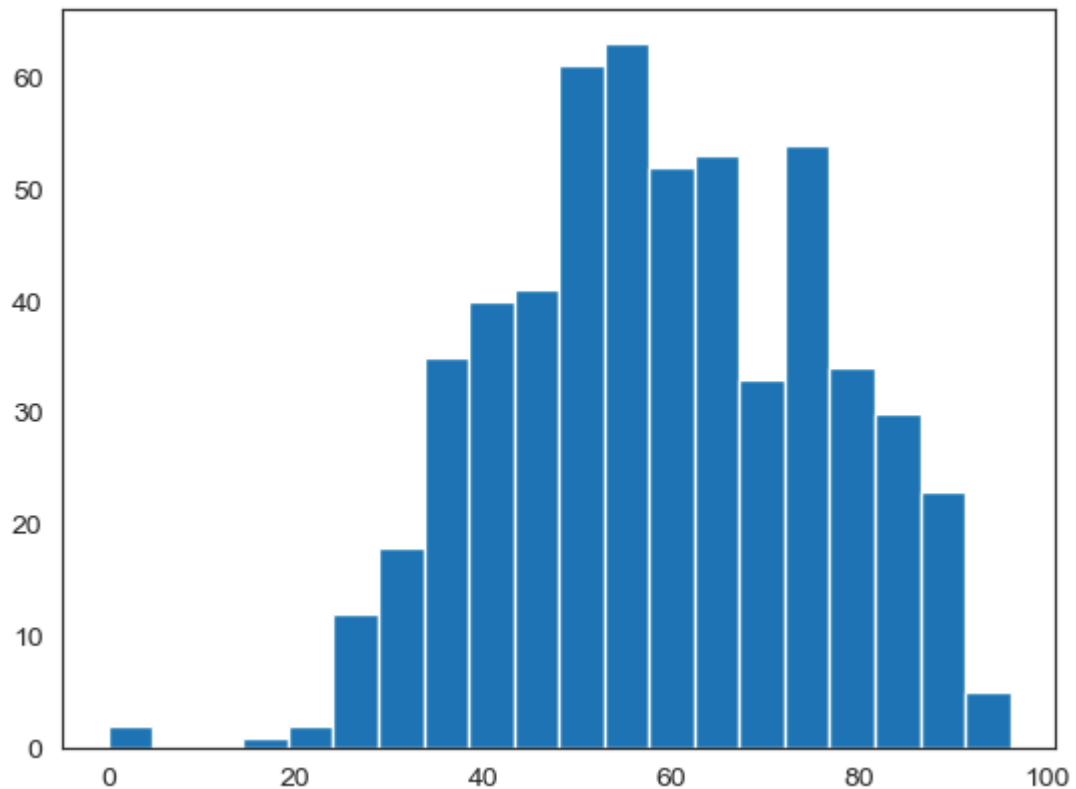
```
In [48]: m2 = sns.distplot(movies.AudienceRating, bins = 15)
```



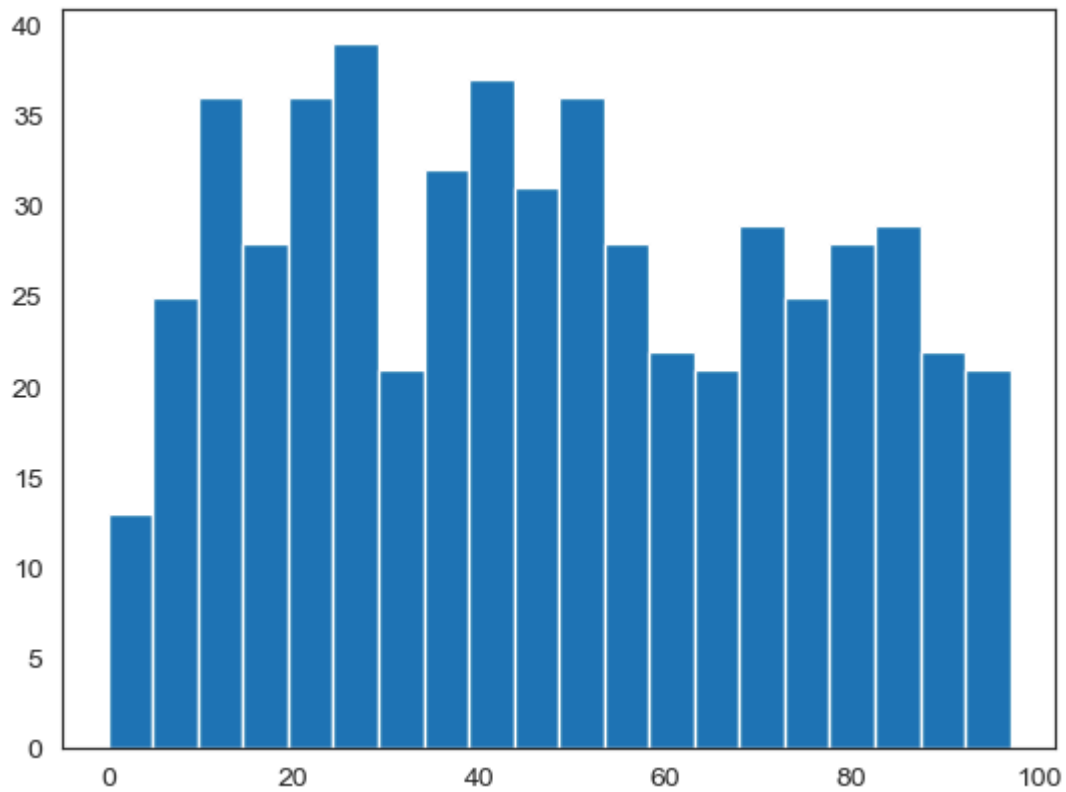
```
In [50]: #sns.set_style('darkgrid')
n1 = plt.hist(movies.AudienceRating, bins=15)
```



```
In [52]: sns.set_style('white') #normal distribution & called as bell curve  
n1 = plt.hist(movies.AudienceRating, bins=20)
```



```
In [54]: n1 = plt.hist(movies.CriticRating, bins=20) #uniform distribution
```

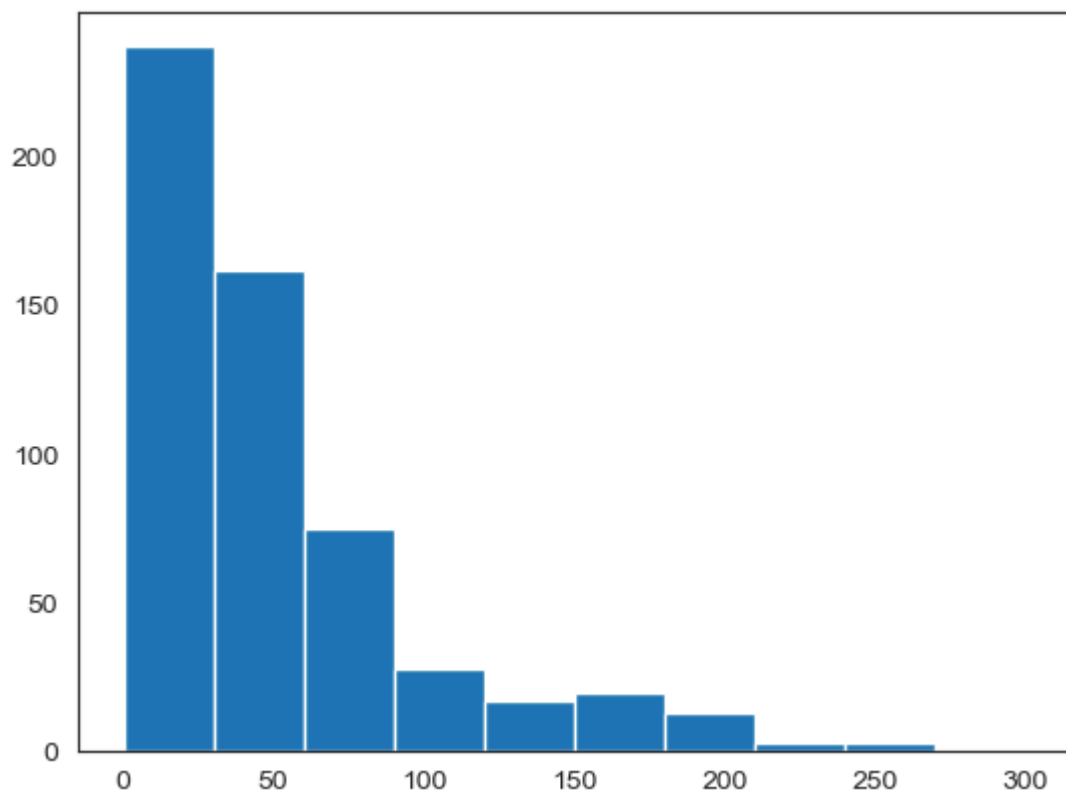


In []: `# <<< chat - 2`

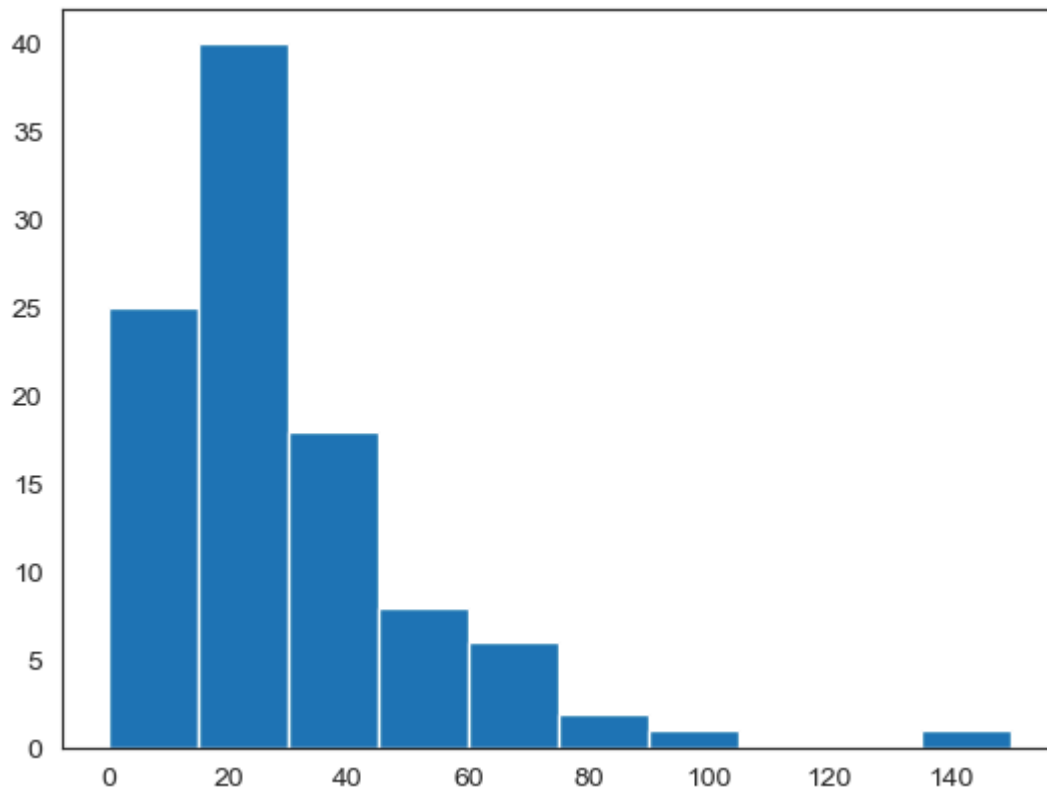
Creating stacked histograms & this is bit tough to understand

In [56]: `#h1 = plt.hist(movies.BudgetMillions)`

```
plt.hist(movies.BudgetMillions)
plt.show()
```



```
In [58]: plt.hist(movies[movies.Genre == 'Drama'].BudgetMillions)
plt.show()
```



```
In [60]: movies.head()
```

```
Out[60]:
```

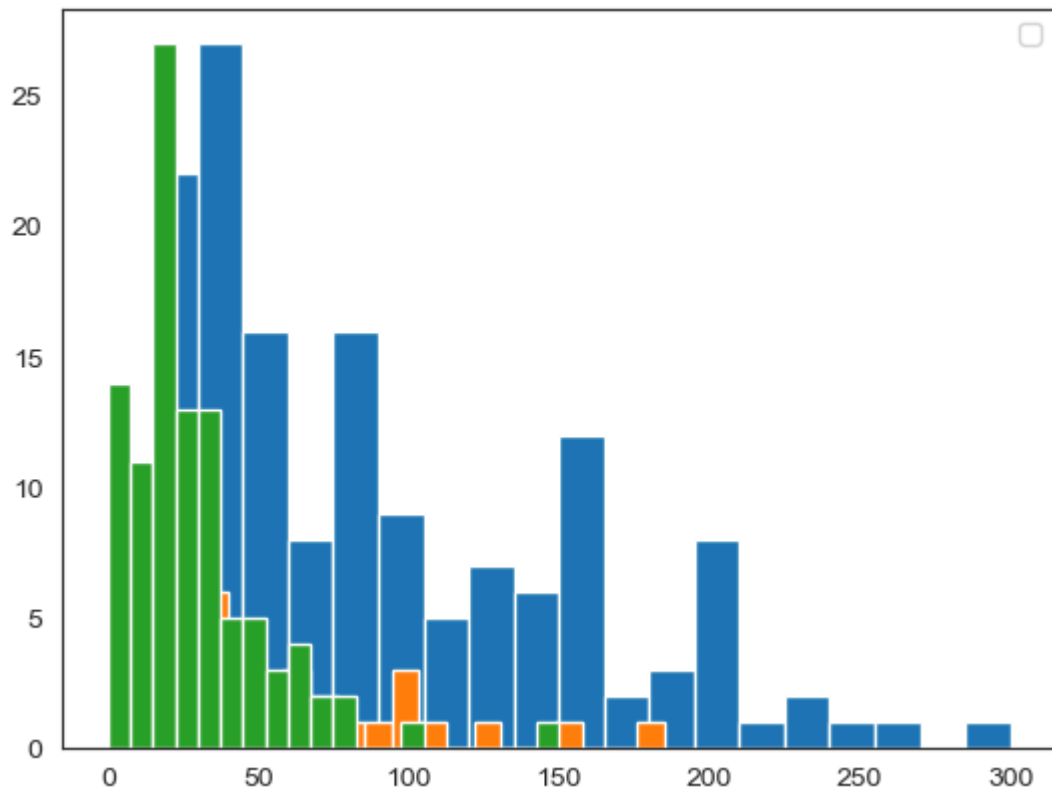
	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
0	(500) Days of Summer	Comedy	87	81	8	2009
1	10,000 B.C.	Adventure	9	44	105	2008
2	12 Rounds	Action	30	52	20	2009
3	127 Hours	Adventure	93	84	18	2010
4	17 Again	Comedy	55	70	20	2009

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

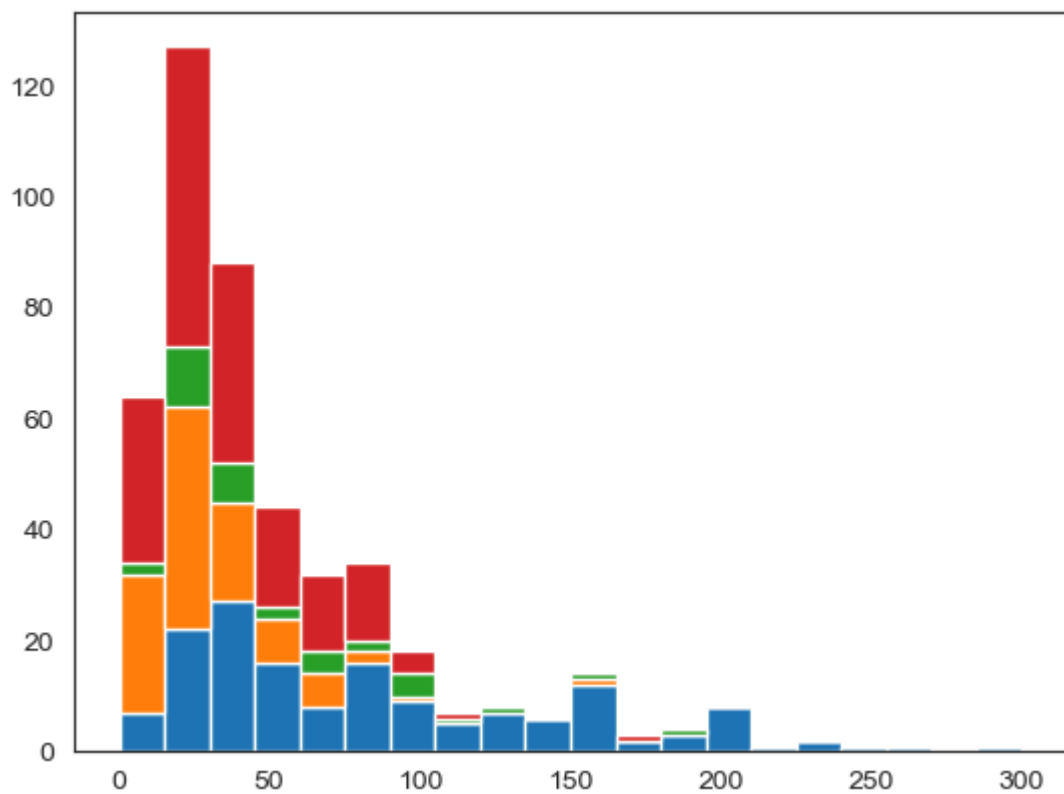
```
In [64]: # Below plots are stacked histogram becuae overlaped

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()
```

No artists with labels found to put in legend. Note that artists whose label start with an underscore are ignored when legend() is called with no argument.



```
In [66]: plt.hist([movies[movies.Genre == 'Action'].BudgetMillions, \
    movies[movies.Genre == 'Drama'].BudgetMillions, \
    movies[movies.Genre == 'Thriller'].BudgetMillions, \
    movies[movies.Genre == 'Comedy'].BudgetMillions],
    bins = 20, stacked = True)
plt.show()
```

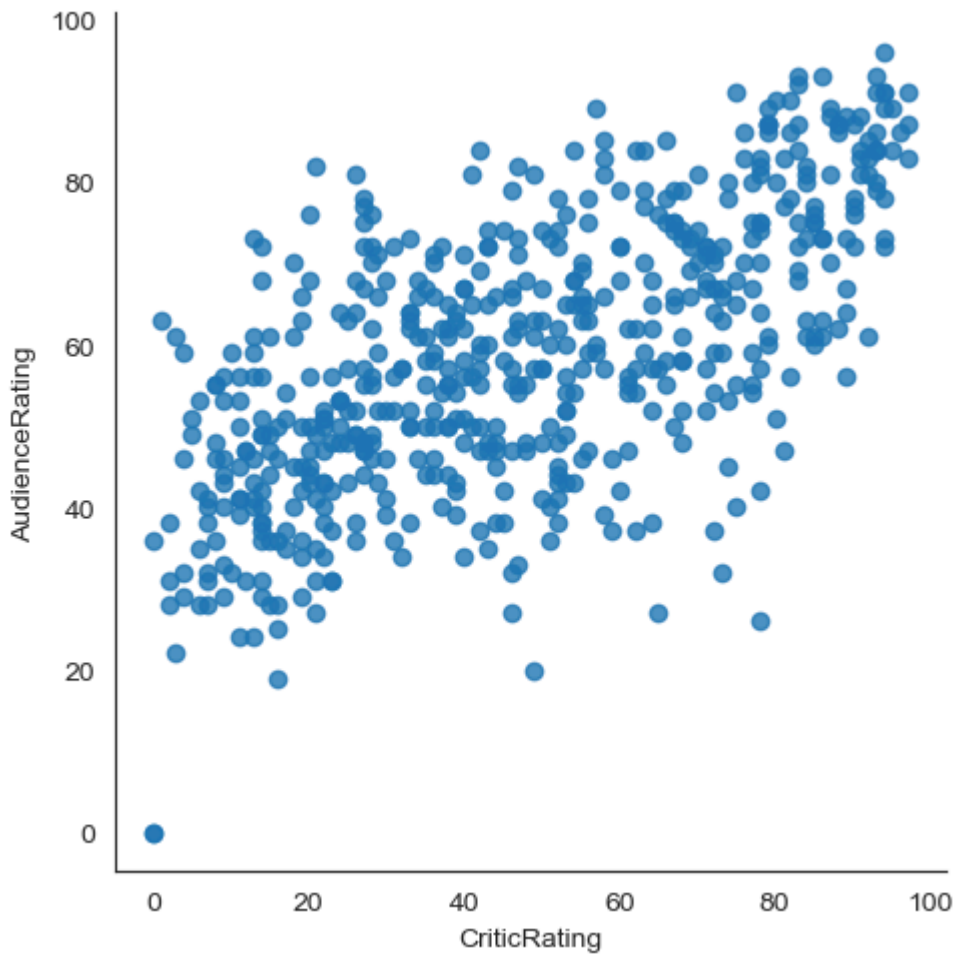


```
In [68]: # if you have 100 categories you cannot copy & paste all the things
```

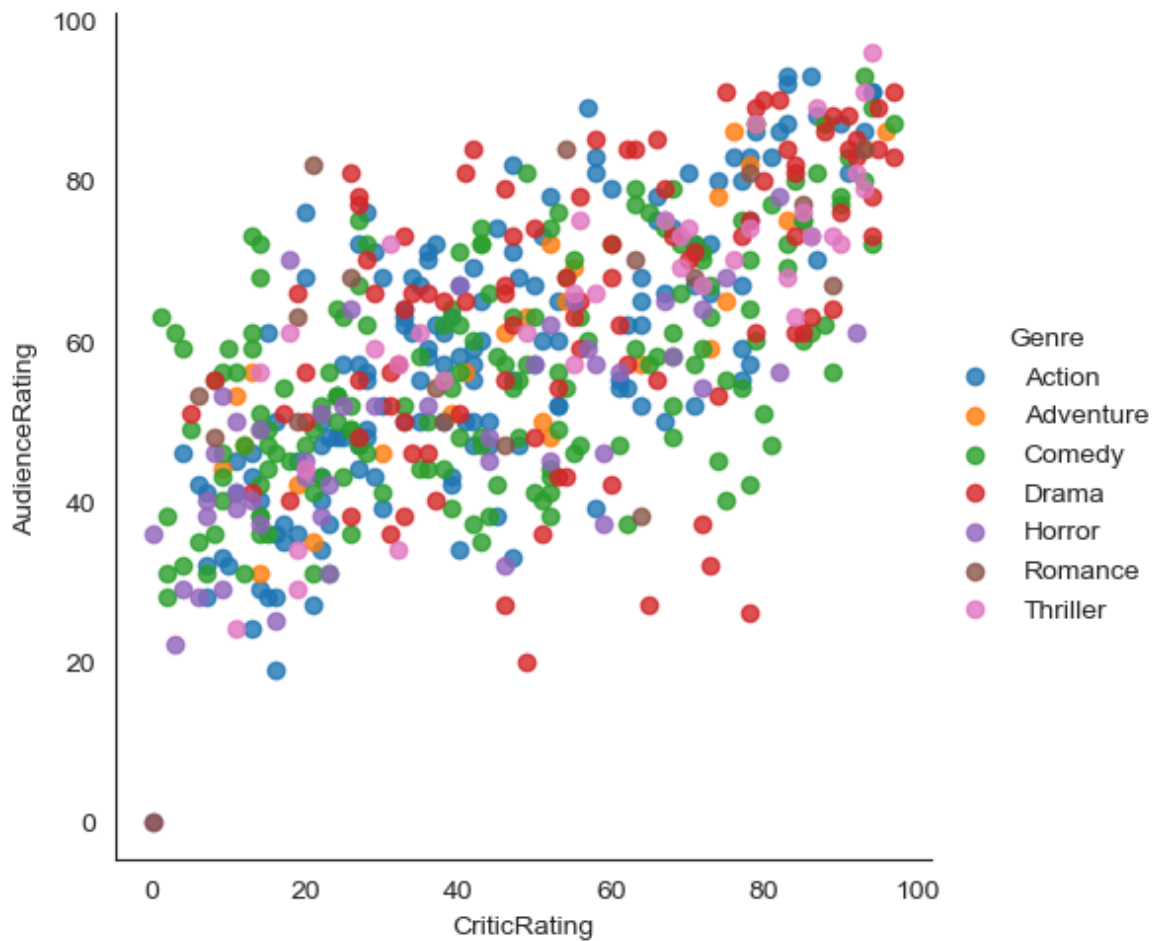
```
for gen in movies.Genre.cat.categories:  
    print(gen)
```

Action
Adventure
Comedy
Drama
Horror
Romance
Thriller

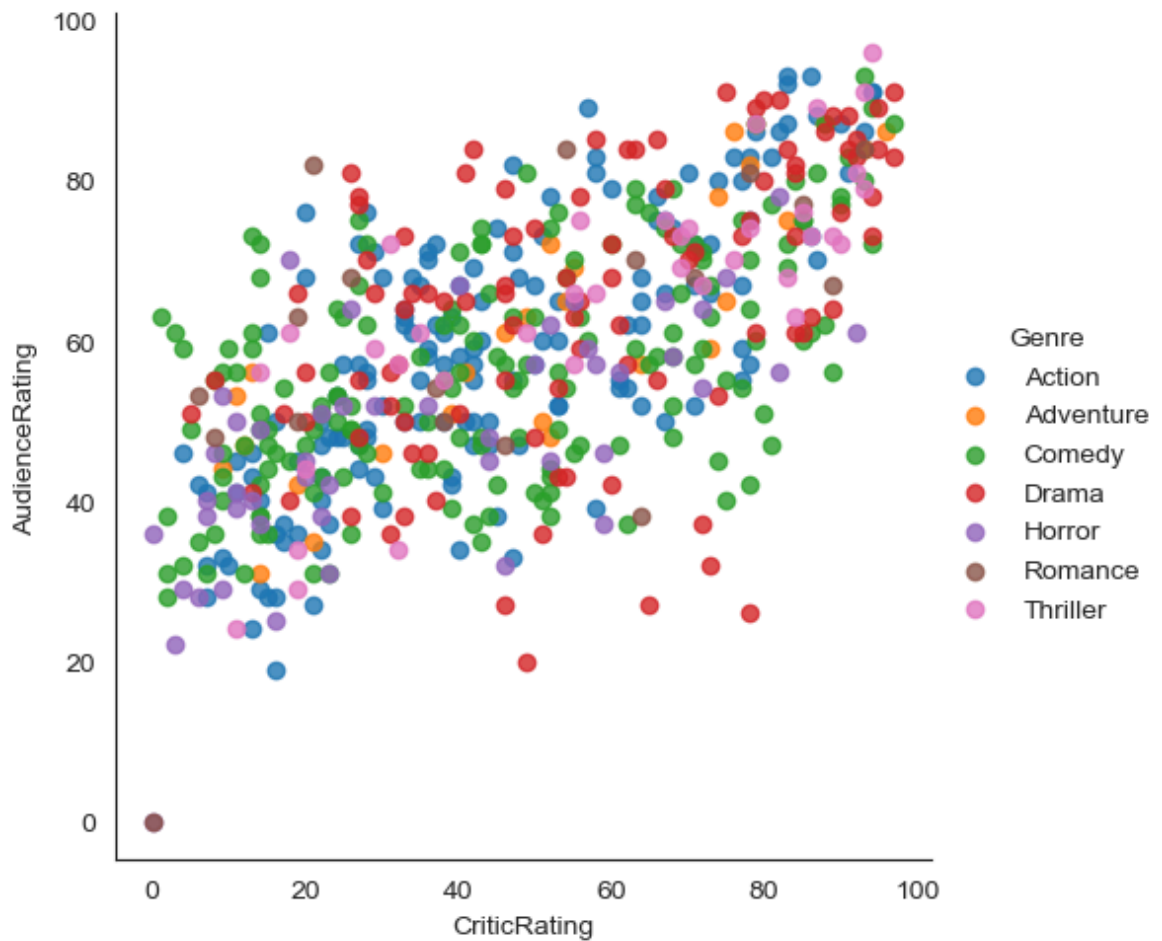
```
In [70]: vis1 = sns.lmplot(data=movies, x='CriticRating', y='AudienceRating',\  
                           fit_reg=False)
```



```
In [72]: vis1 = sns.lmplot(data=movies, x='CriticRating', y='AudienceRating',\  
                           fit_reg=False, hue = 'Genre')
```



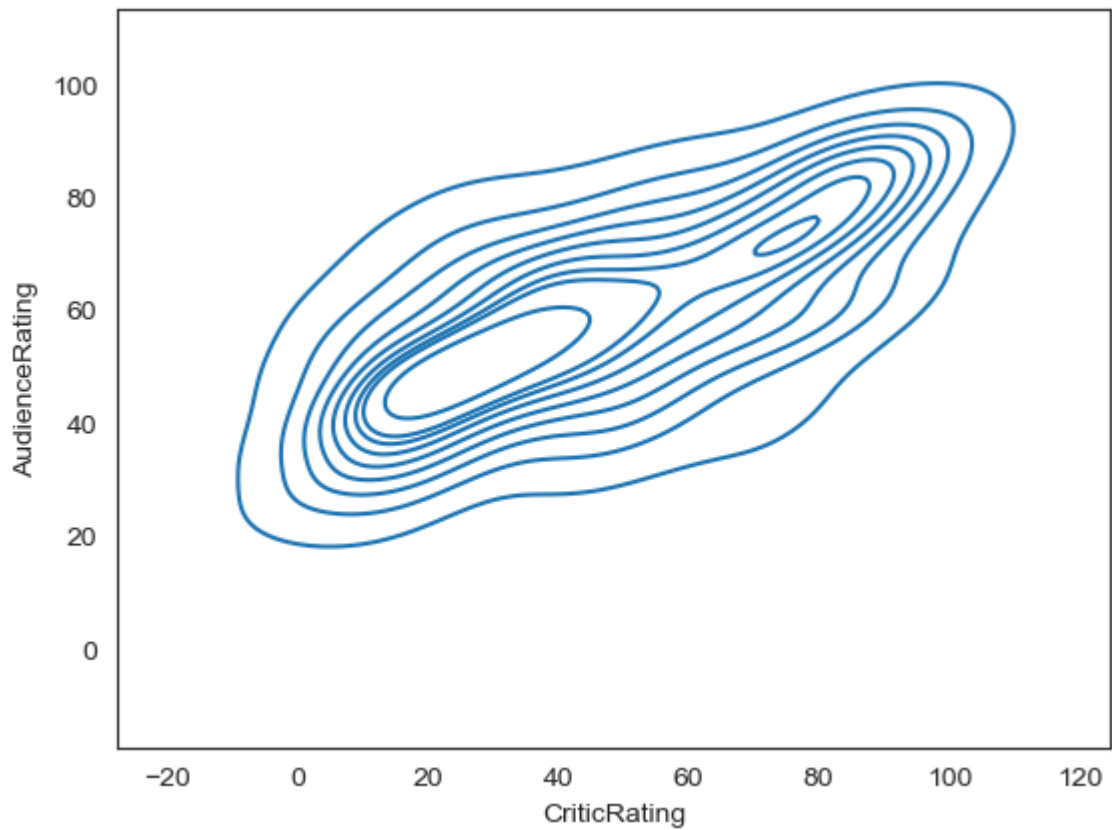
```
In [74]: vis = sns.lmplot(data=movies, x='CriticRating', y='AudienceRating',\
                        fit_reg=False, hue = 'Genre', aspect=1)
```



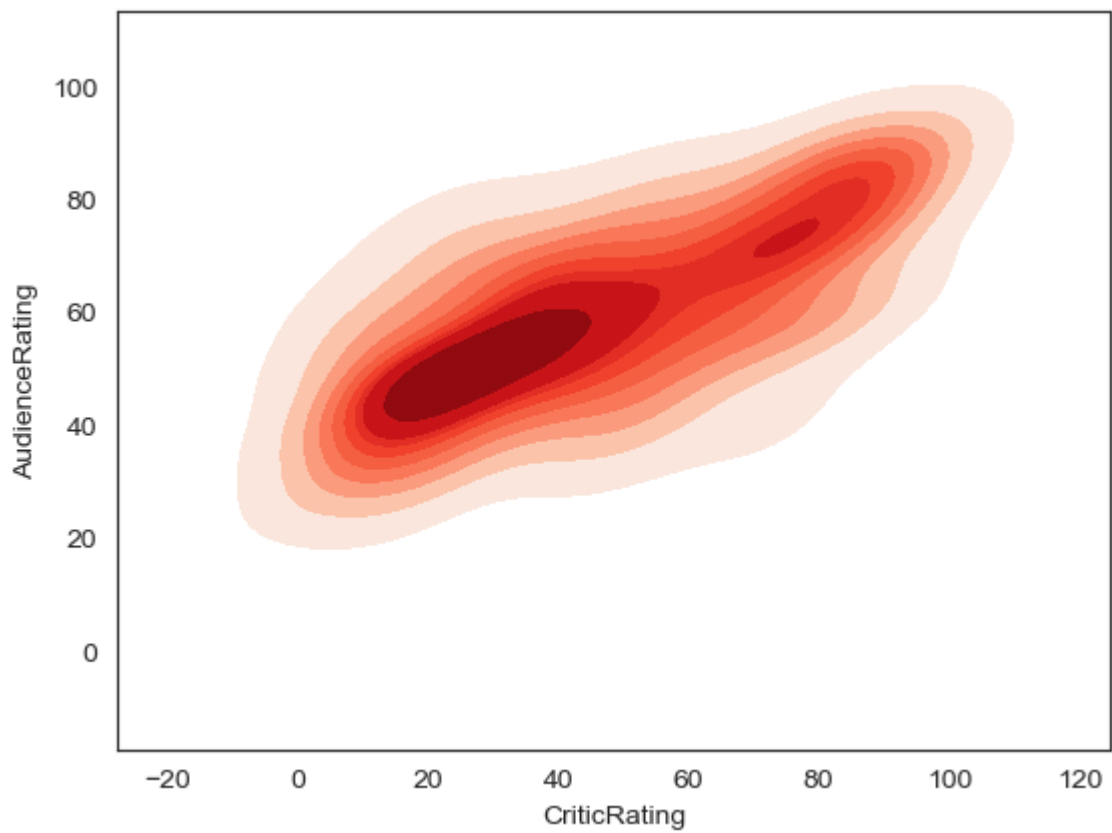
```
In [ ]: # Kernal Density Estimate plot ( KDE PLOT)
# how can i visulize audience rating & critics rating . using scatterplot
```

```
In [76]: k1 = sns.kdeplot(x = movies.CriticRating,y= movies.AudienceRating)

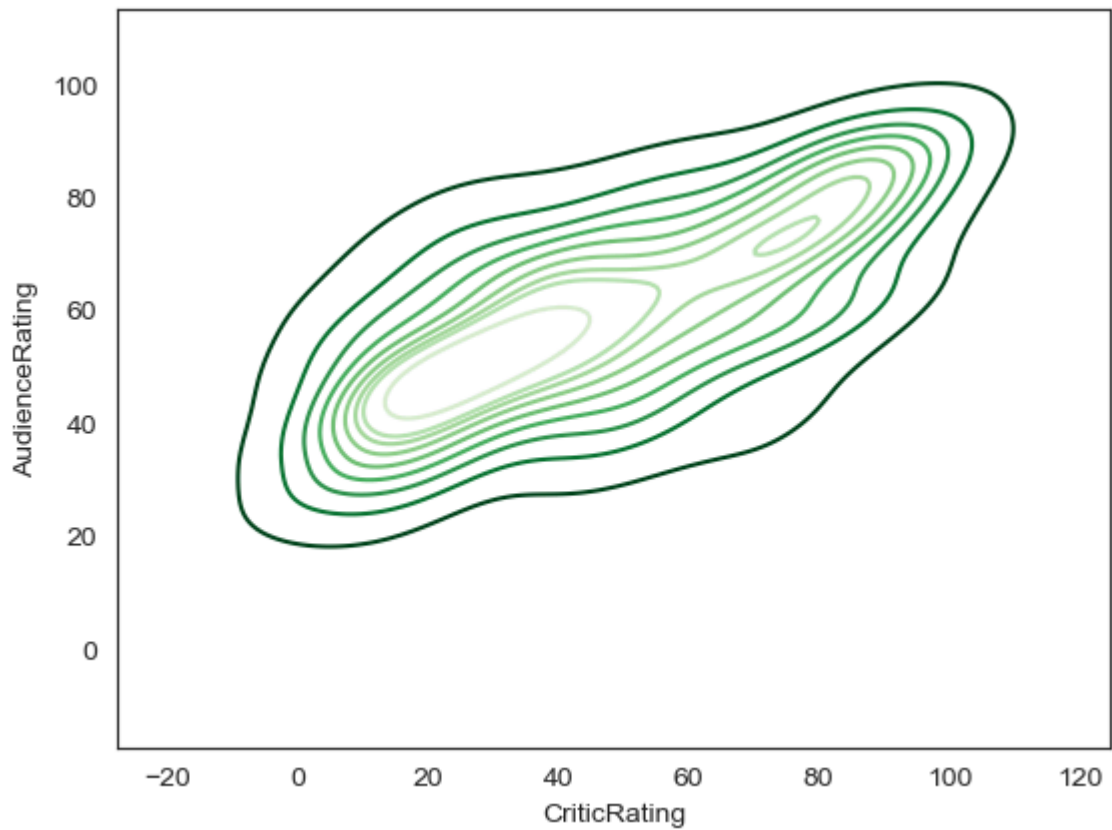
# where do u find more density and how density is distibuted across from the the c
# center point is kernal this is calld KDE & insteade of dots it visualize like th
# we can able to clearly see the spread at the audience ratings
```

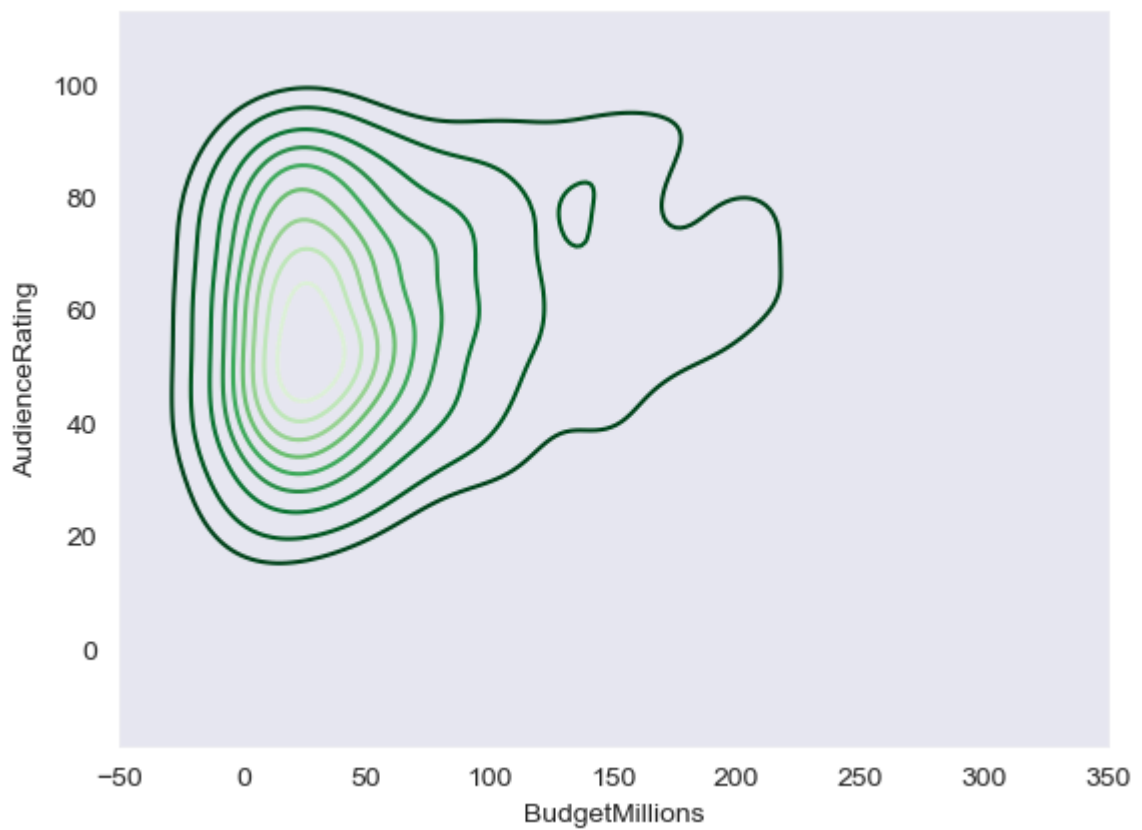
```
In [78]: k1 = sns.kdeplot(x=movies.CriticRating,y=movies.AudienceRating,shade = True,shade_
```



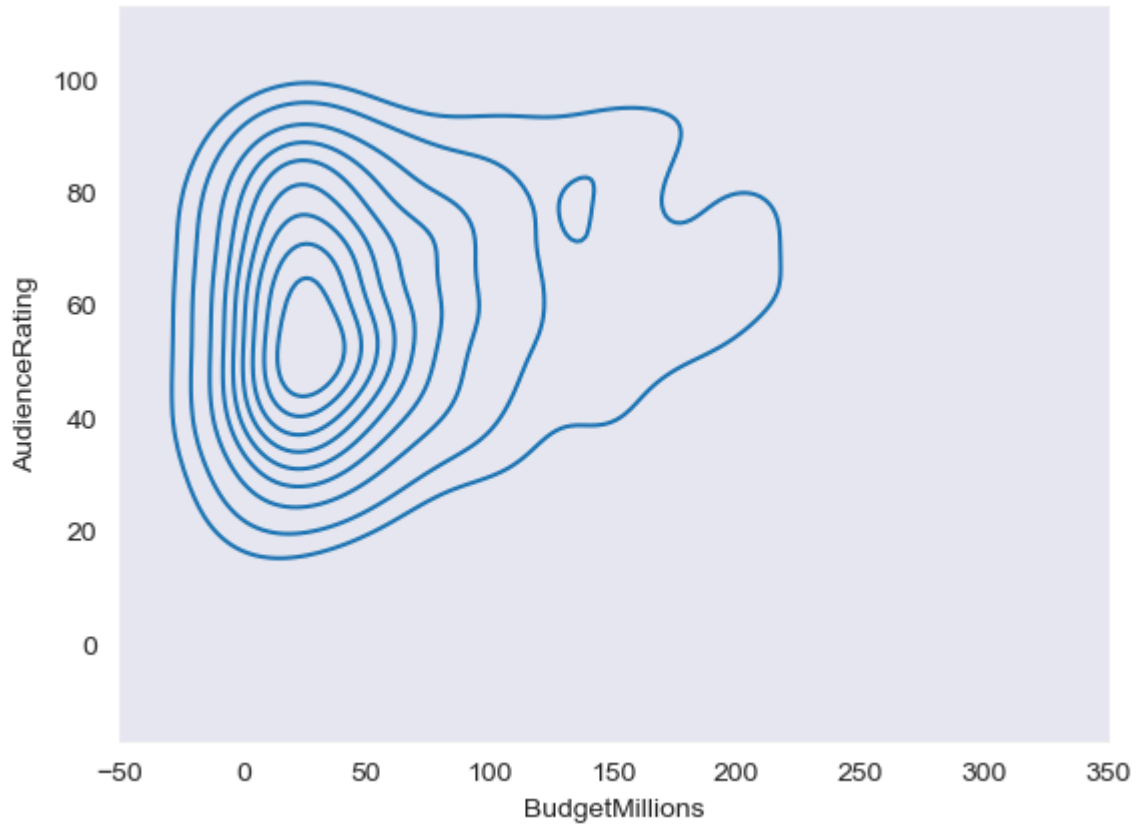
```
In [80]: k2 = sns.kdeplot(x=movies.CriticRating,y=movies.AudienceRating,shade_lowest=False,
```



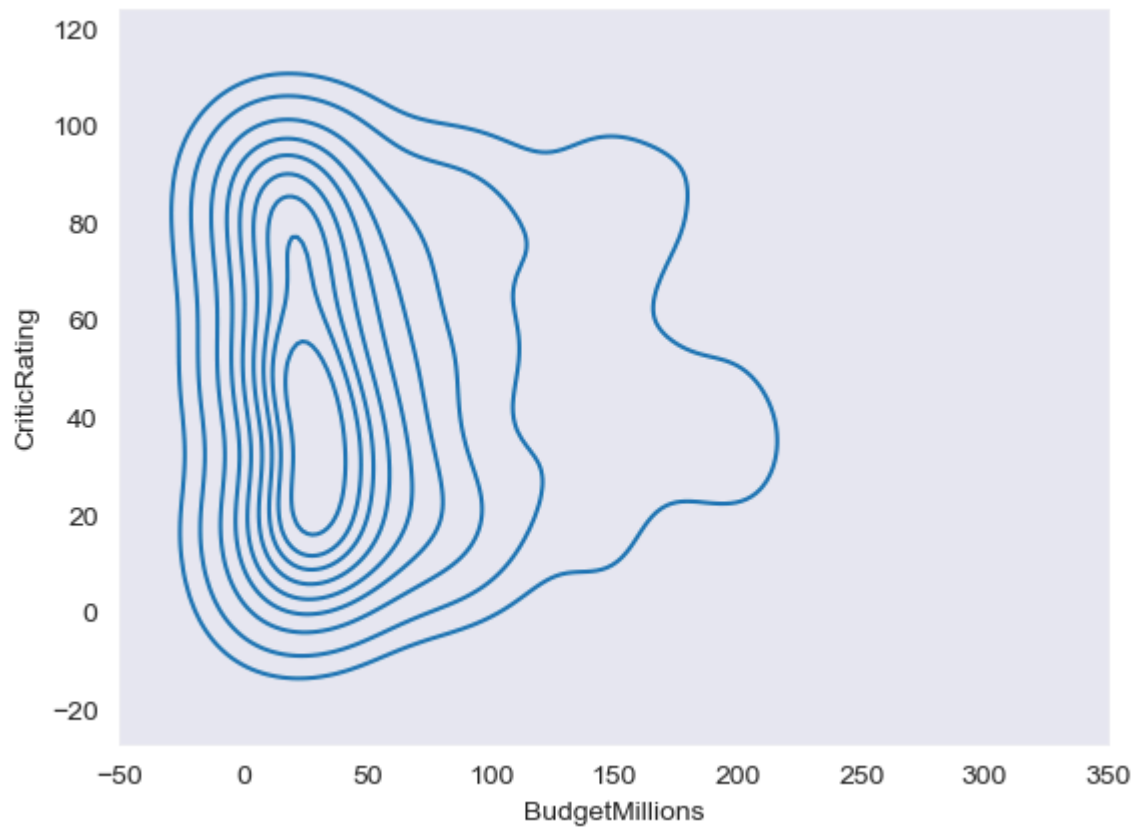
```
In [82]: sns.set_style('dark')
k1 = sns.kdeplot(x=movies.BudgetMillions,y=movies.AudienceRating,shade_lowest=False)
```



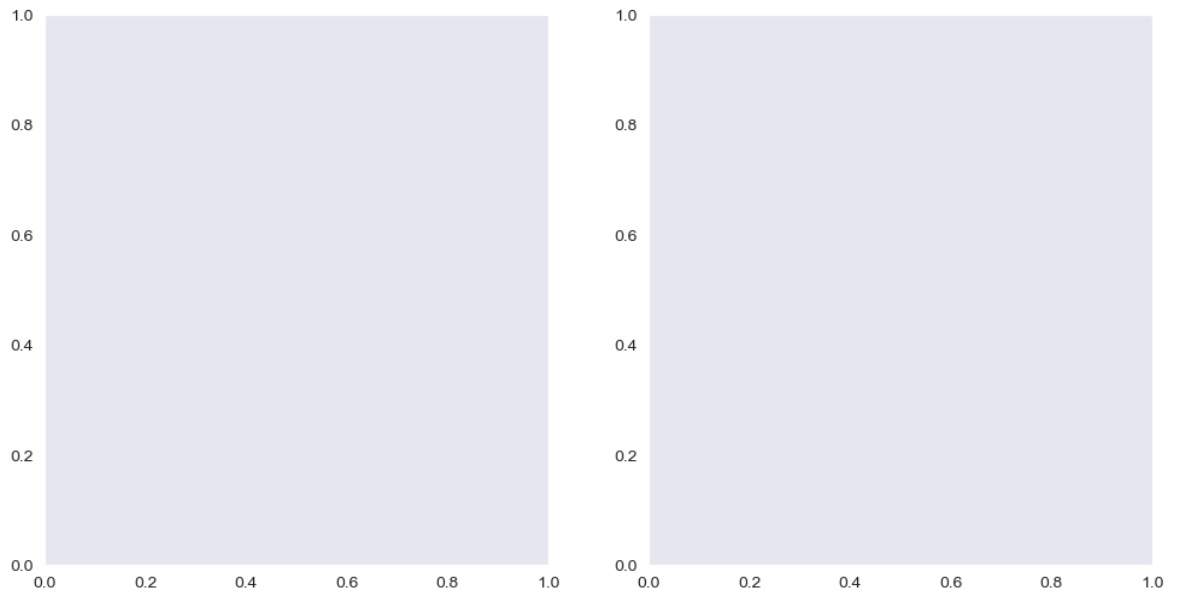
```
In [84]: sns.set_style('dark')
k1 = sns.kdeplot(x=movies.BudgetMillions,y=movies.AudienceRating)
```



```
In [86]: k2 = sns.kdeplot(x=movies.BudgetMillions,y=movies.CriticRating)
```

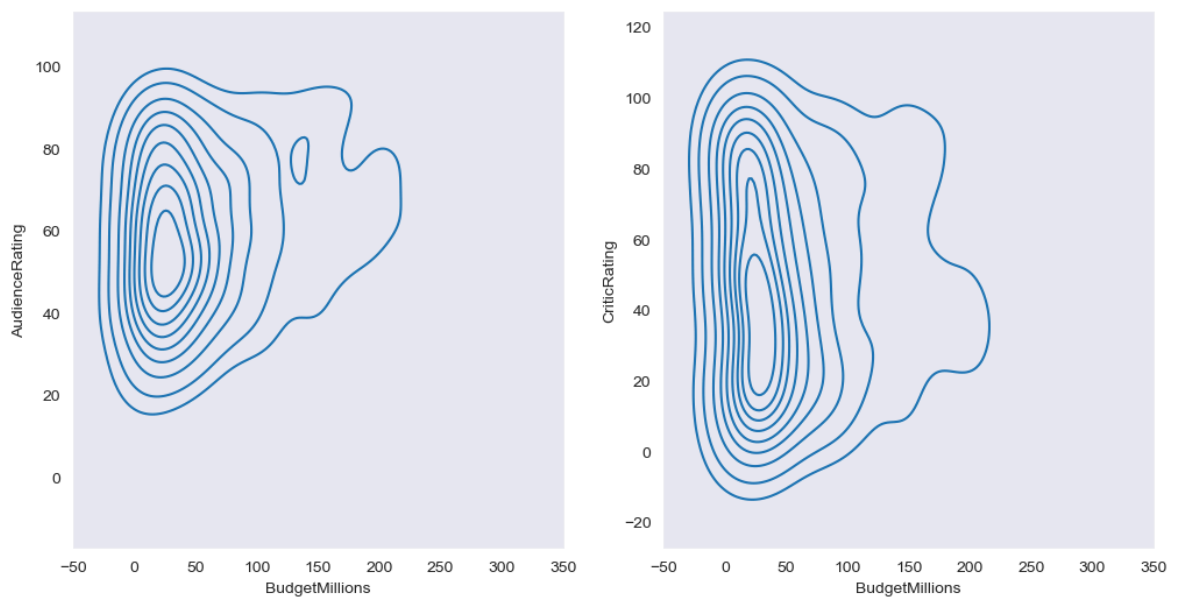


```
In [88]: #subplots  
  
f, ax = plt.subplots(1,2, figsize =(12,6))  
#f, ax = plt.subplots(3,3, figsize =(12,6))
```



```
In [90]: f, axes = plt.subplots(1,2, figsize =(12,6))

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

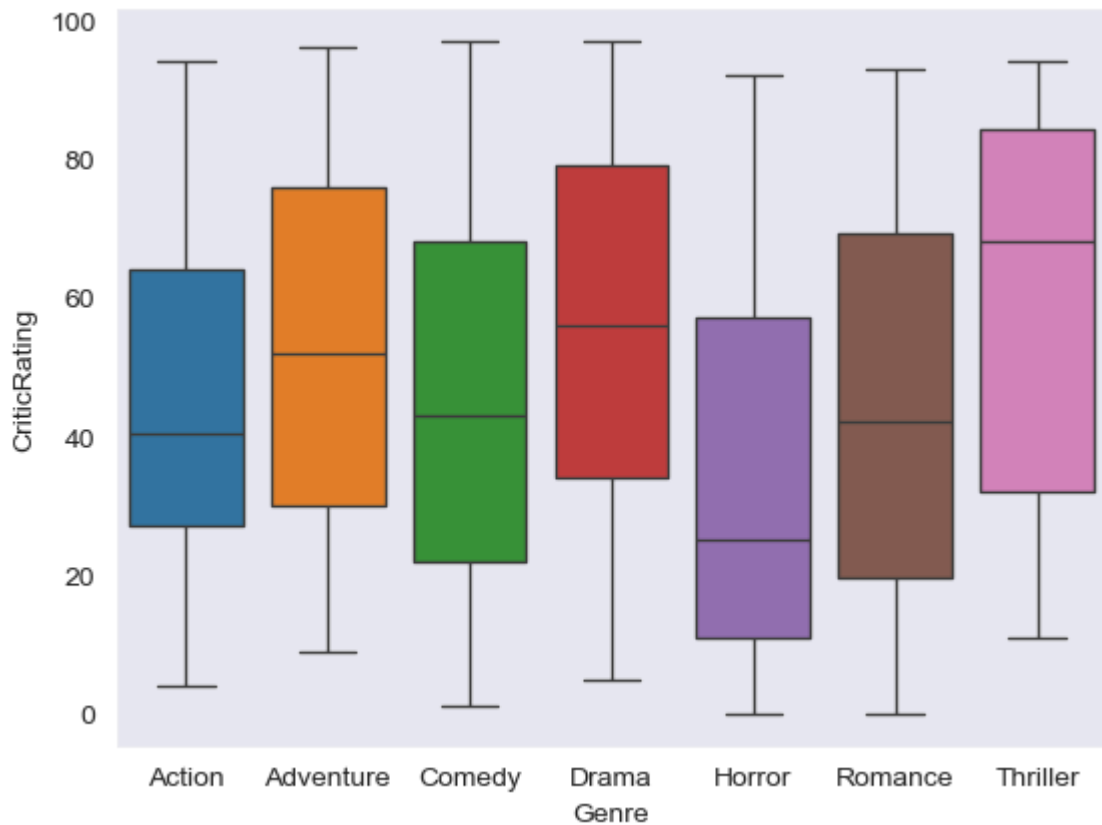


```
In [92]: axes
```

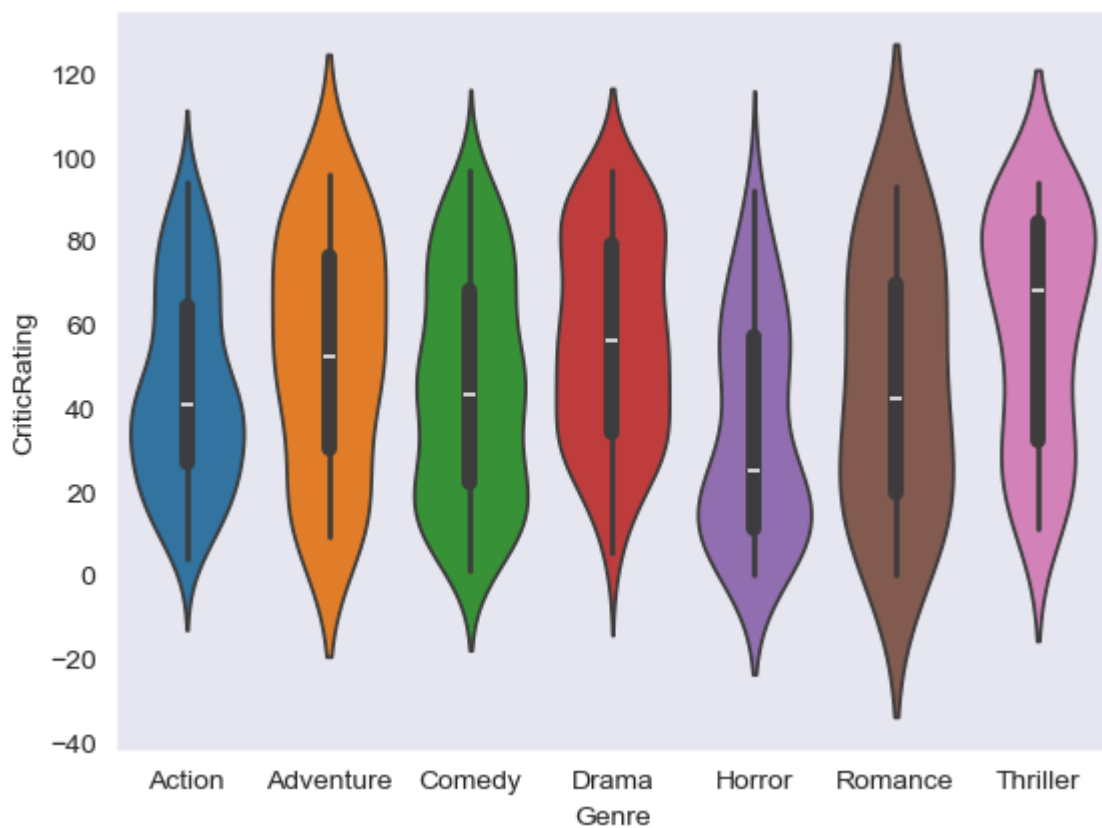
```
Out[92]: array([<Axes: xlabel='BudgetMillions', ylabel='AudienceRating'>,
                <Axes: xlabel='BudgetMillions', ylabel='CriticRating'>],
              dtype=object)
```

```
In [94]: #Box plots -

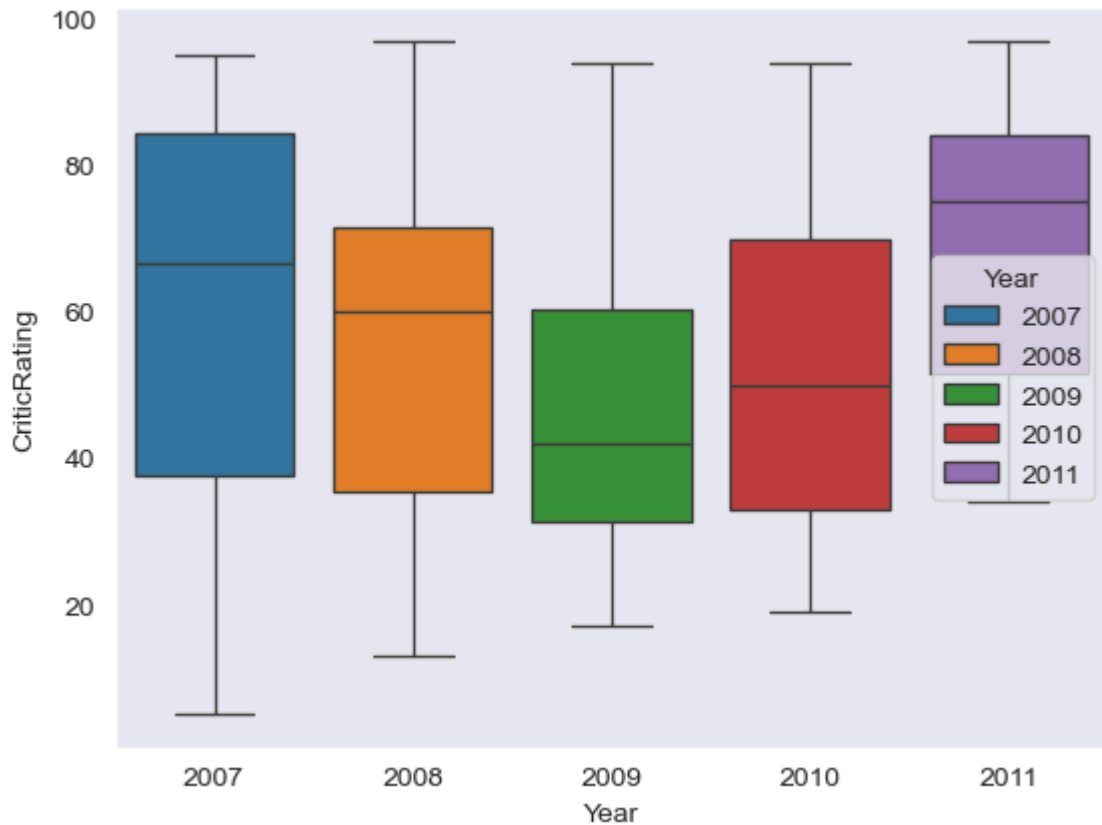
w = sns.boxplot(data=movies, x='Genre', y = 'CriticRating',hue='Genre')
```



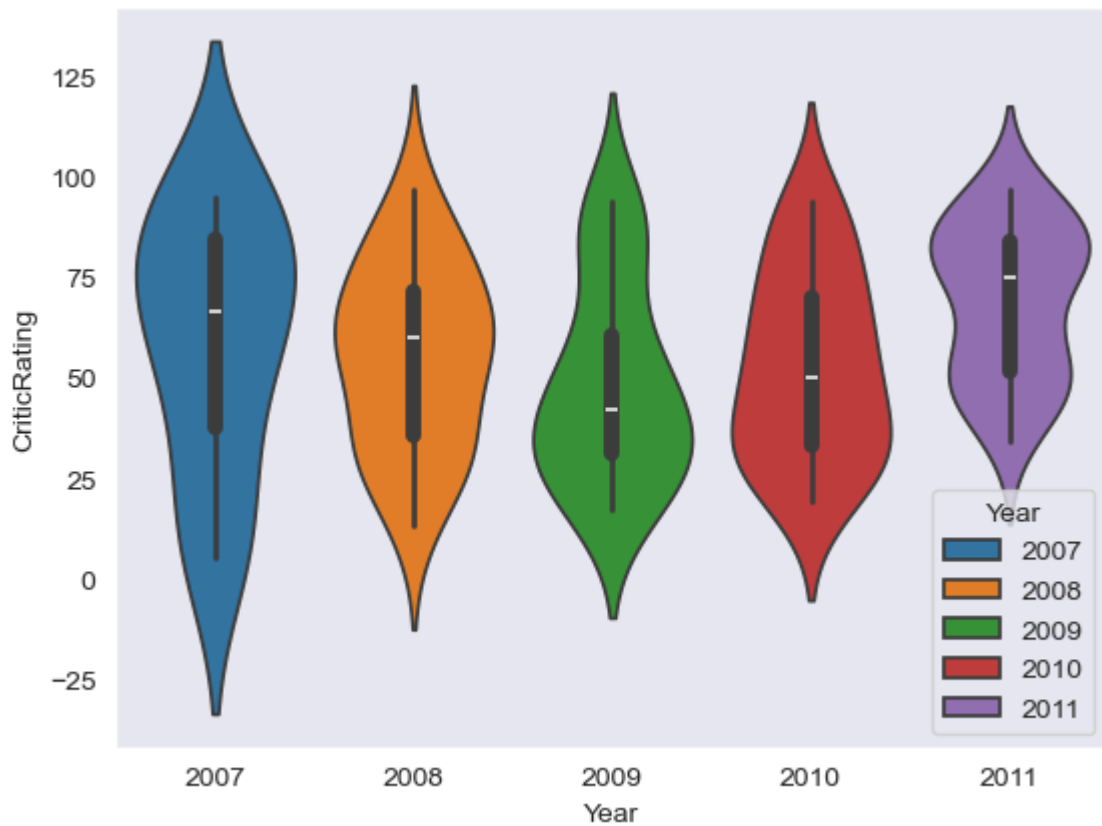
```
In [96]: z = sns.violinplot(data=movies, x='Genre', y='CriticRating', hue='Genre')  
plt.show()
```



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

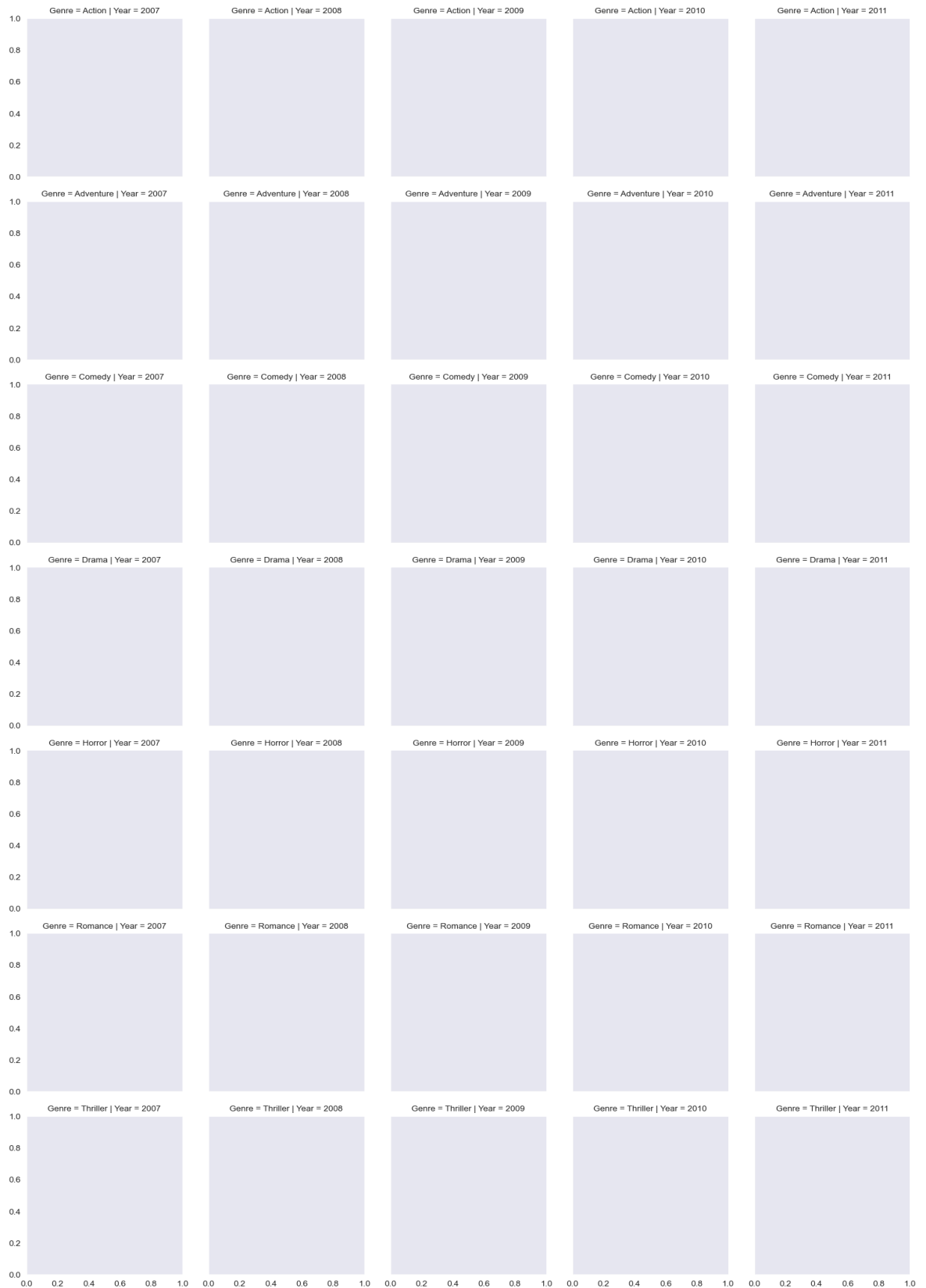


```
In [100]: z = sns.violinplot(data=movies[movies.Genre == 'Drama'], x='Year', y = 'CriticRating')
```



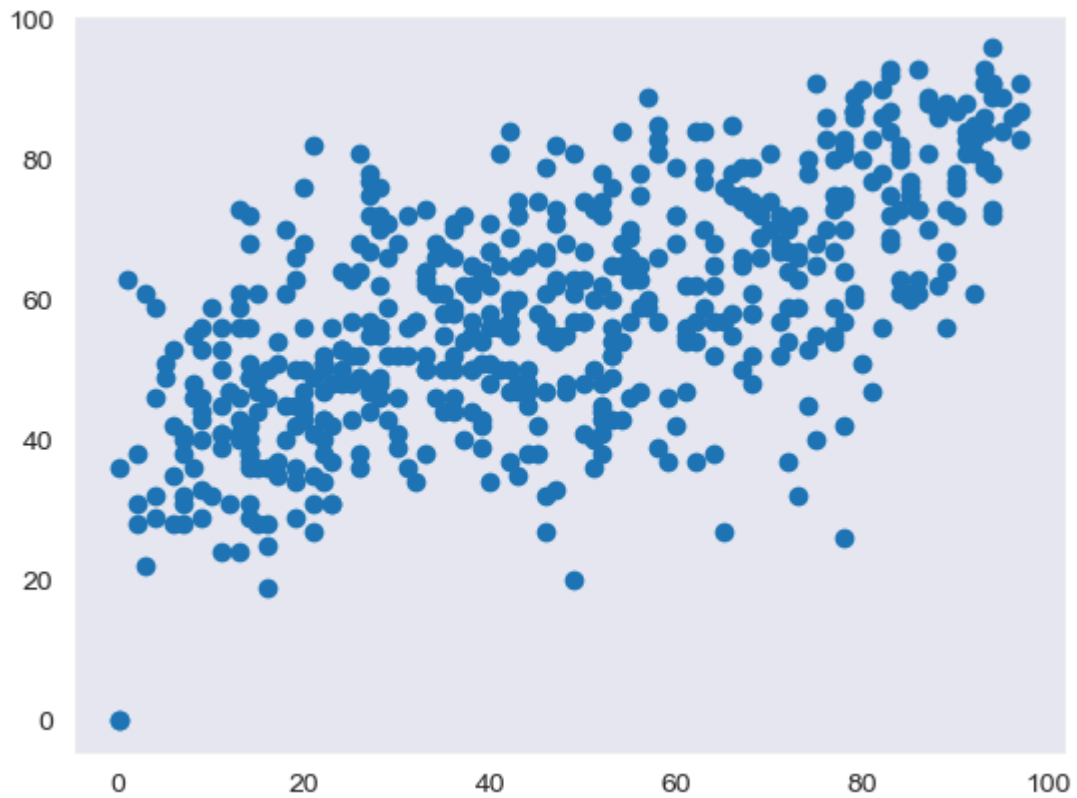
```
In [ ]: # Createing a Facet grid
```

```
In [104]: g =sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre') #kind of subplot
```

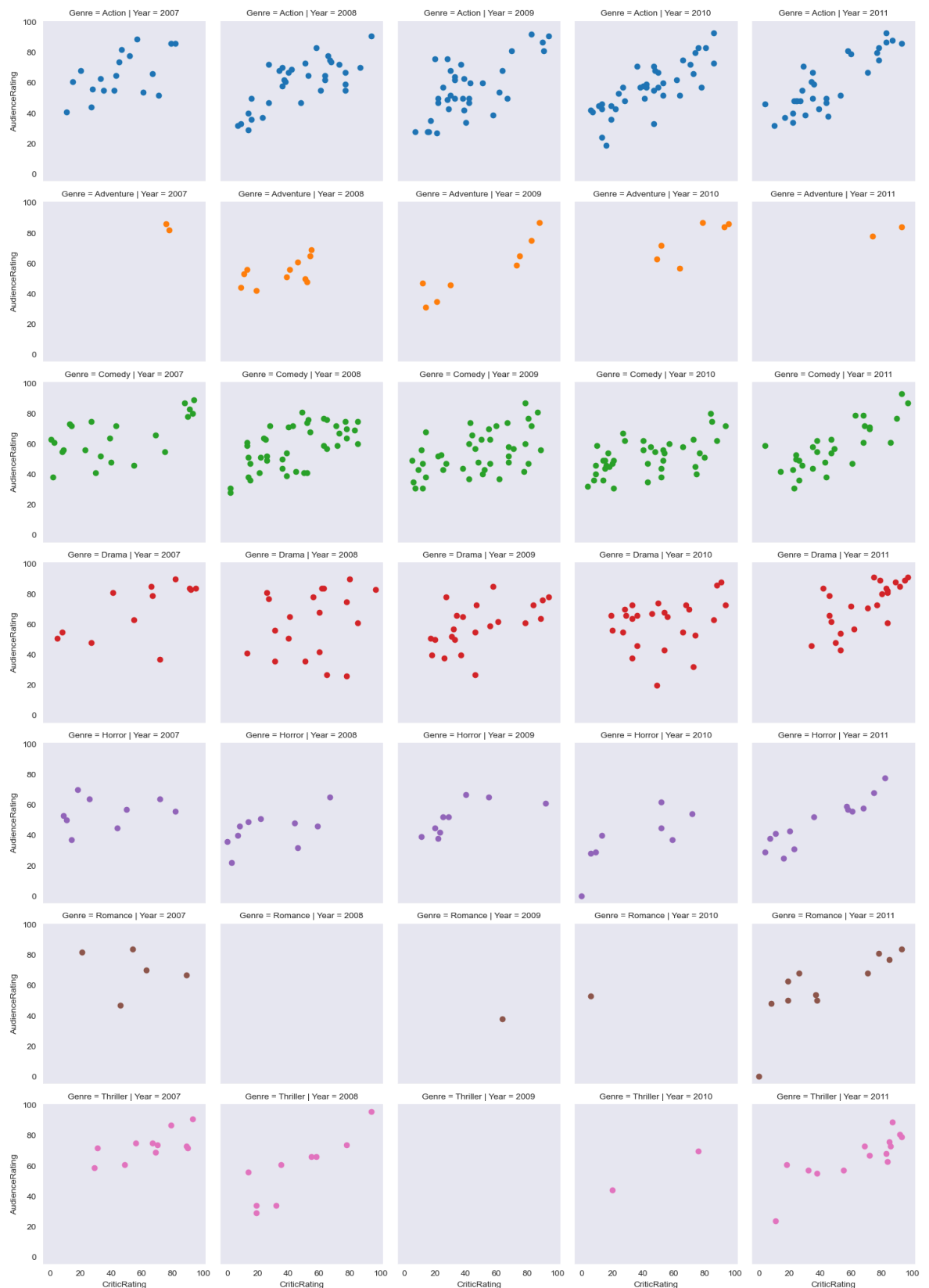


```
In [108]: plt.scatter(x=movies.CriticRating,y=movies.AudienceRating)
```

```
Out[108]: <matplotlib.collections.PathCollection at 0x18c4bd36090>
```



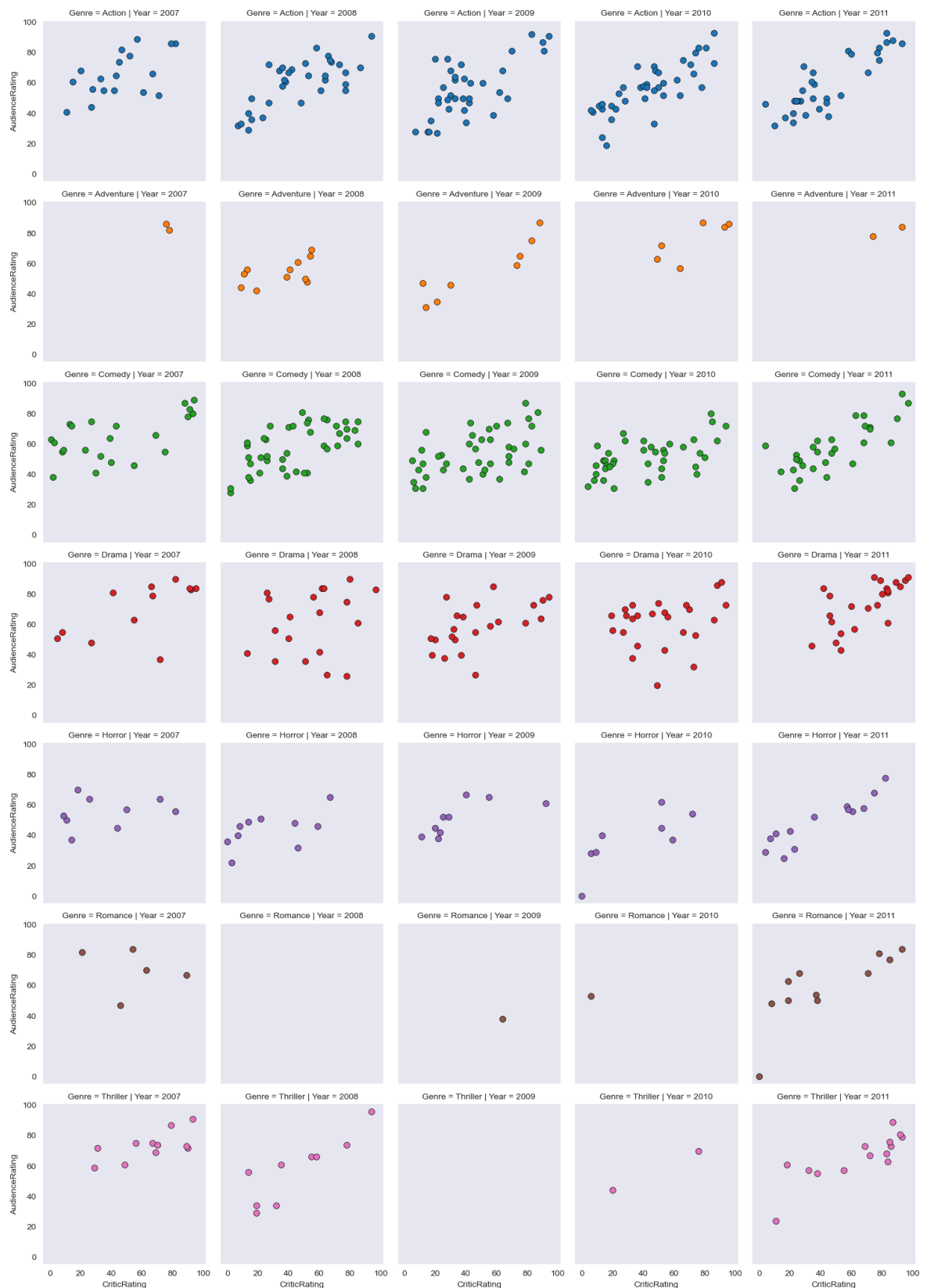
```
In [110]: g = sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre')
g = g.map(plt.scatter, 'CriticRating', 'AudienceRating' )
#scatterplots are mapped in facetgrid
```

```
In [112]: # 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
```



```
In [114]: #
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
```



```
In [138]: # python is not vectorize programming language
# Building dashboards (dashboard - combination of chats)

sns.set_style('darkgrid')
f, axes = plt.subplots (2,2, figsize = (15,15))

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

k1.set(xlim=(-20,160))
```

```

k2.set(xlim=(-20,160))

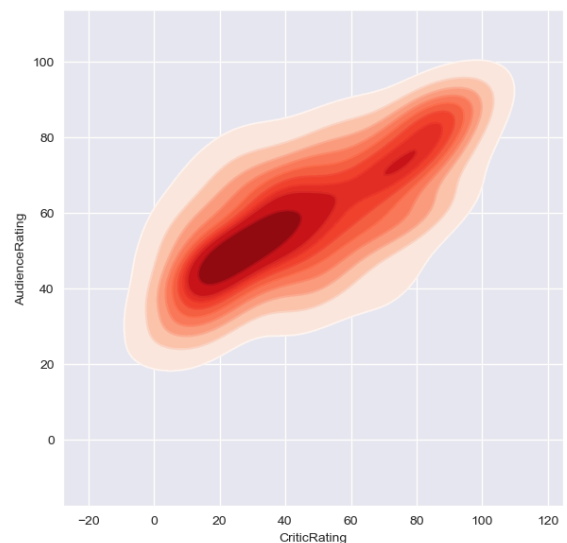
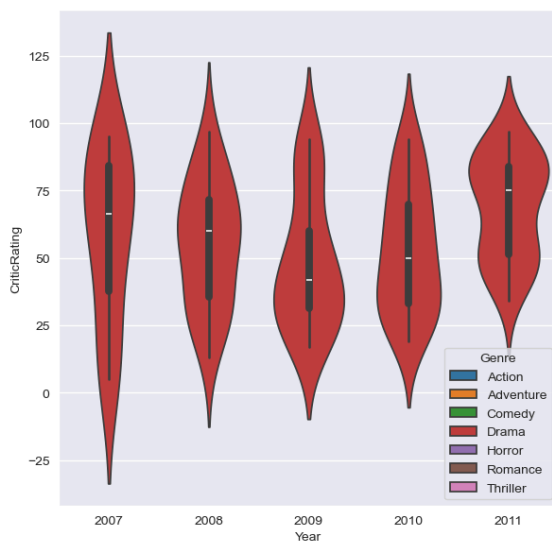
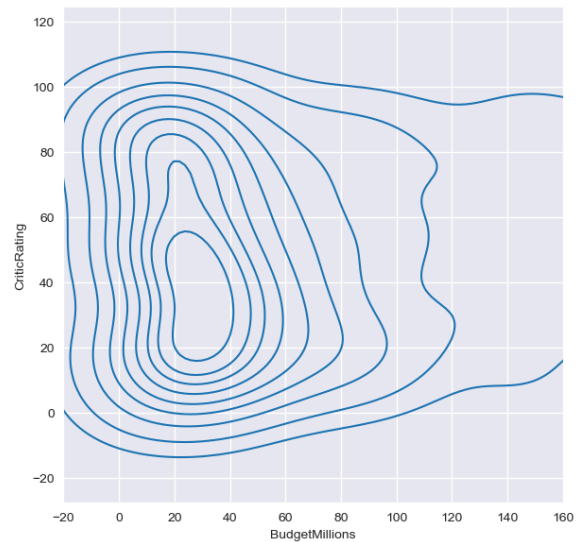
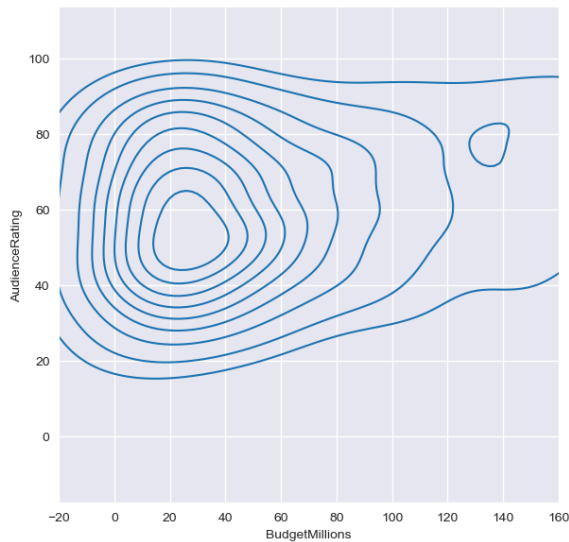
z = sns.violinplot(data=movies[movies.Genre=='Drama'], x='Year', y = 'CriticRating')

k4 = sns.kdeplot(x=movies.CriticRating,y=movies.AudienceRating,shade = True,shade_

k4b = sns.kdeplot(x=movies.CriticRating,y=movies.AudienceRating,cmap='Reds',ax = a

plt.show()

```

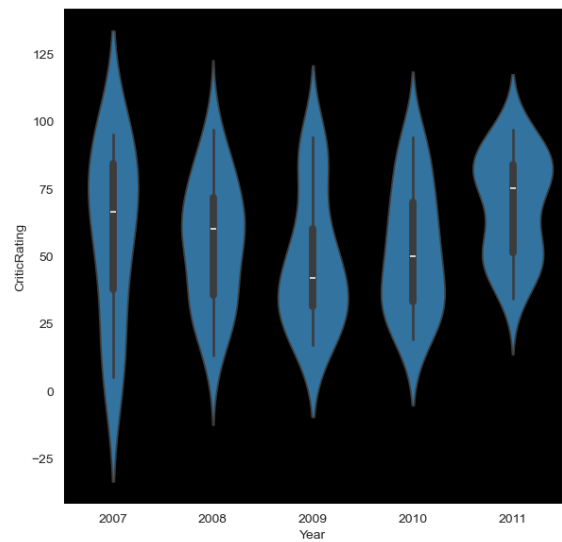
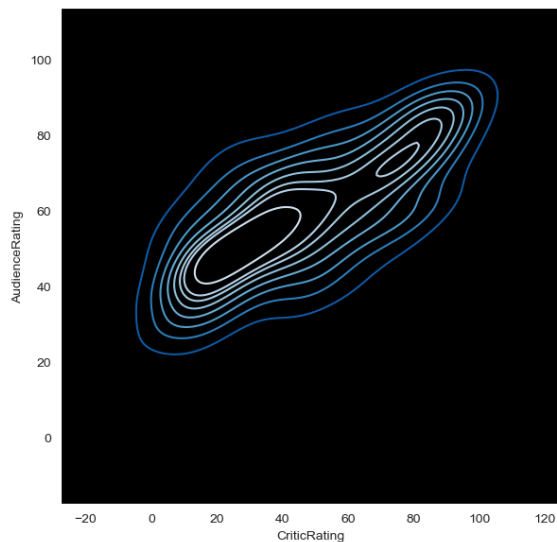
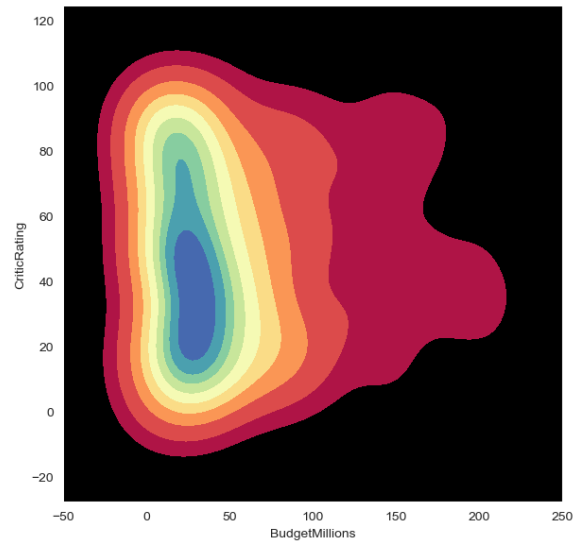
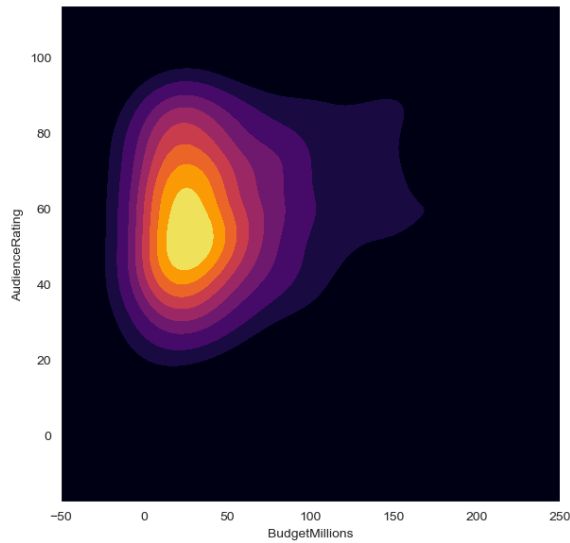


```

In [144]: sns.set_style('dark',{'axes.facecolor':'black'})
f, axes = plt.subplots(2,2,figsize=(15,15))
#plot [0,0]
k1 = sns.kdeplot(x=movies.BudgetMillions,y=movies.AudienceRating, \
shade = True, shade_lowest=True,cmap = 'inferno', \
ax = axes[0,0])
#plot [0,1]
k2 = sns.kdeplot(x=movies.BudgetMillions, y=movies.CriticRating,\
shade=True, shade_lowest=False, cmap='Spectral',\
ax = axes[0,1])
#plot[1,0]
k3 = sns.kdeplot(x=movies.CriticRating, y=movies.AudienceRating, \
shade = False,shade_lowest=True,cmap='Blues_r', \
ax=axes[1,0])
#plot[1,1]

```

```
vi = sns.violinplot(data=movies[movies.Genre=='Drama'], \
x='Year', y = 'CriticRating', ax=axes[1,1])
k1.set(xlim=(-50,250))
k2.set(xlim=(-50,250))
plt.show()
```



In []: 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> Faceted grid
- 9> Building dashboards

In []: *# EDA is completed*