

In [70]: # MOVIE RATING ANALYTICS (ADVANCED VISULATION)

```
import pandas as pd
import os
```

In [71]: os.getcwd() # if you want to change the working directory

Out[71]: 'C:\\\\Users\\\\dell'

In [72]: movies = pd.read_csv(r"C:\\23rd,24th\\MOVIE RATINGS _ ADVANCE VISUALIZATION _ EDA 1.csv")

In [73]: movies

Out[73]:

	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 [74]: len(movies)

Out[74]: 559

In [75]: `movies.head()`

Out[75]:

	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 [76]: `movies.tail()`

Out[76]:

	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 [77]: `movies.columns`

Out[77]: `Index(['Film', 'Genre', 'Rotten Tomatoes Ratings %', 'Audience Ratings %', 'Budget (million $)', 'Year of release'], dtype='object')`

In [79]: `movies.columns = ['Film', 'Genre', 'CriticRating', 'AudienceRating', 'BudgetMillion$']`

In [80]: `movies.head() # Removed spaces & % removed noise characters`

Out[80]:

	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 [81]: `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 [82]: `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[82]:

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

Out[83]:

```
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 [84]: `movies.Film`

Out[84]:

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

In [86]: `movies.Film`

Out[86]:

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 [87]: `movies.head()`

Out[87]:

	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 [88]: `movies.info()`

```
# now the same thing we will change genres 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 [89]: `movies.Genre = movies.Genre.astype('category')`
`movies.Year = movies.Year.astype('category')`

In [90]: `movies.Genre`

```
Out[90]: 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 [91]: `movies.Year # is it real no. year you can take average,min,max but out come have`

```
Out[91]: 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 [92]: `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 [93]: `movies.Genre.cat.categories`

```
Out[93]: Index(['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance',
               'Thriller'],
               dtype='object')
```

In [94]: `movies.describe()`

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

Out[94]:

	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

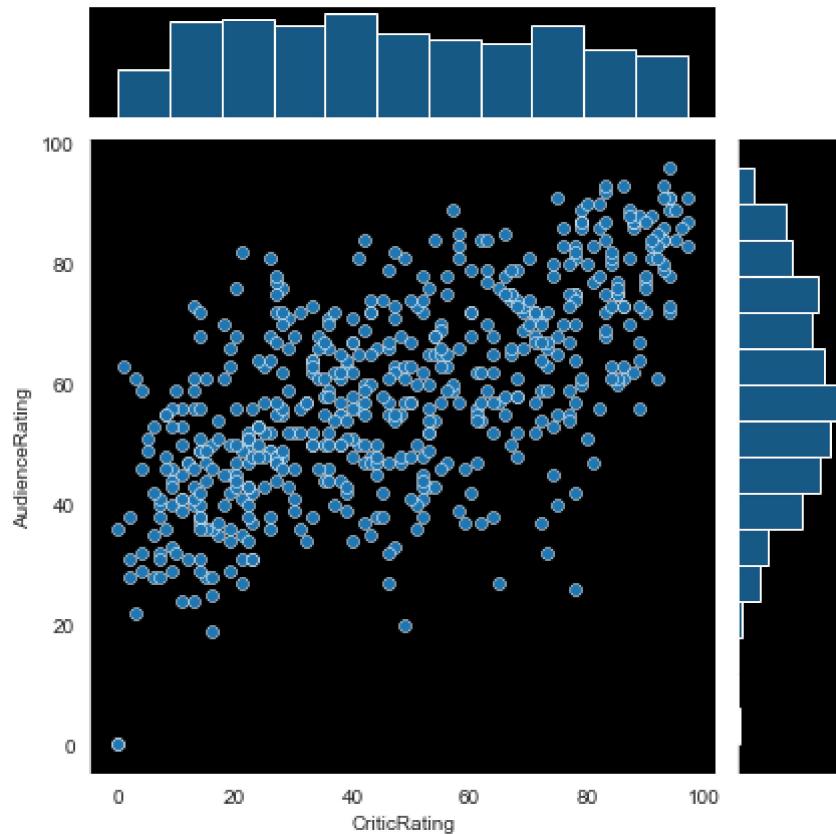
In [95]: *# How to working with joint plots*

```
from matplotlib import pyplot as plt
import seaborn as sns
%matplotlib inline
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)

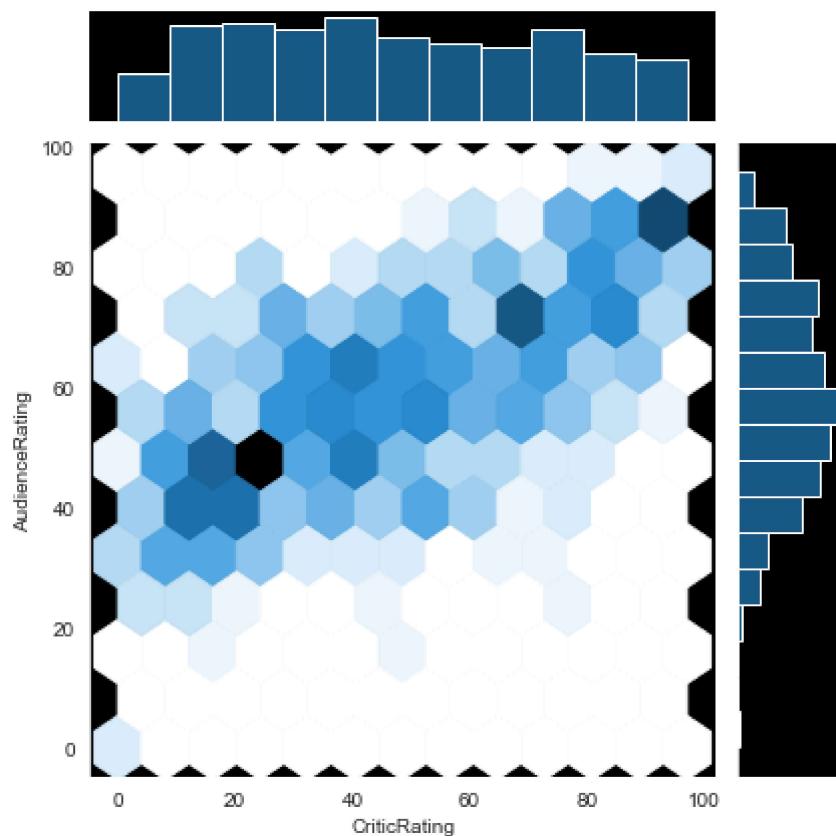
In [96]:

```
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 likelihood to watch audience rating
# Let me explain the excel - if you filter audience rating & critic rating. critics rating
```

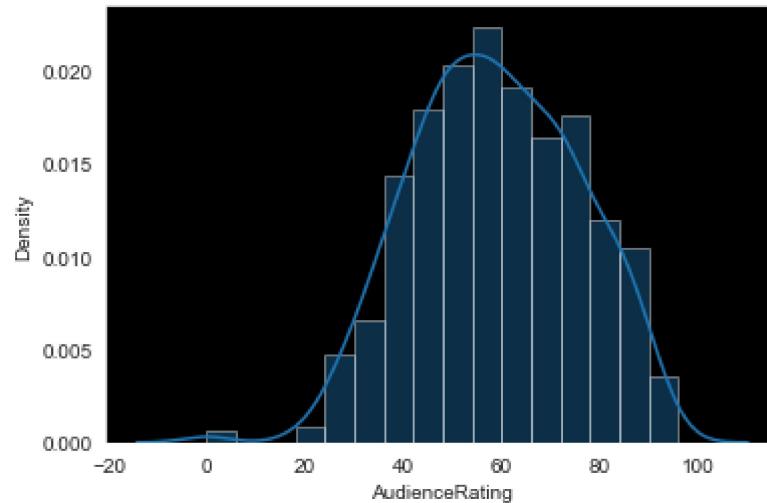


```
In [97]: j = sns.jointplot( data = movies, x = 'CriticRating', y = 'AudienceRating', kind=
```

```
#j = sns.jointplot( data = movies, x = 'CriticRating', y = 'AudienceRating', kind=
```

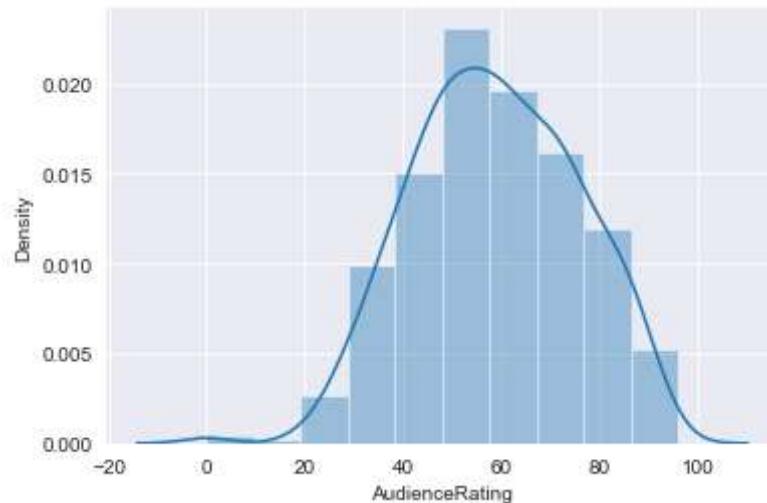


```
In [98]: #Histograms  
  
# <<< chat1  
  
m1 = sns.distplot(movies.AudienceRating)  
  
#y - axis generated by seaborn automatically that is the powefull of seaborn gall
```

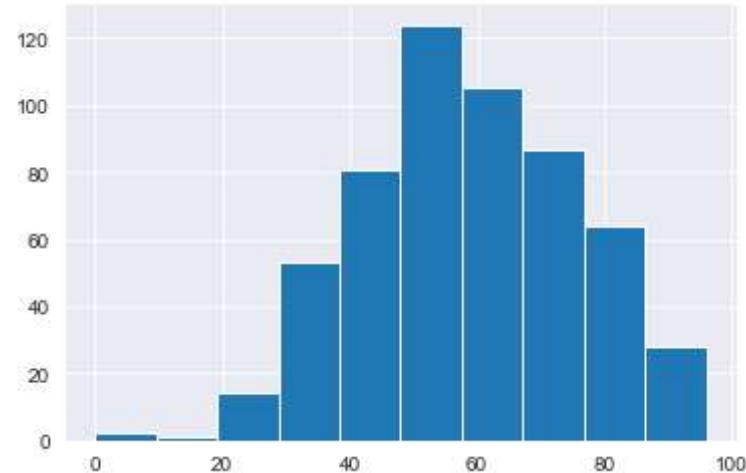


```
In [99]: sns.set_style('darkgrid')
```

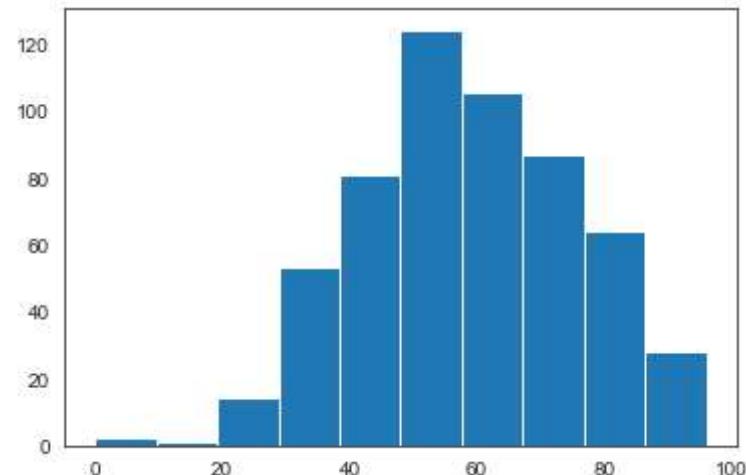
```
In [100]: m2 = sns.distplot(movies.AudienceRating, bins = 10)
```



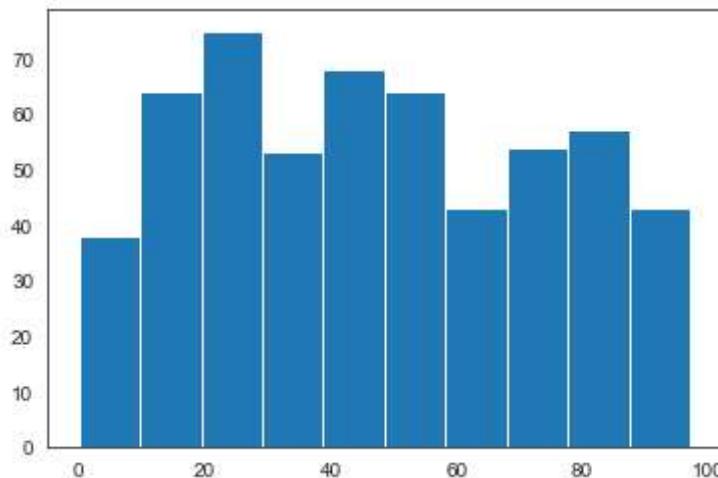
```
In [101]: sns.set_style('darkgrid')
n1 = plt.hist(movies.AudienceRating, bins=10)
```



```
In [102]: sns.set_style('white') #normal distribution & called as bell curve
n1 = plt.hist(movies.AudienceRating, bins=10)
```



```
In [103]: n1 = plt.hist(movies.CriticRating, bins=10) #uniform distribution
```

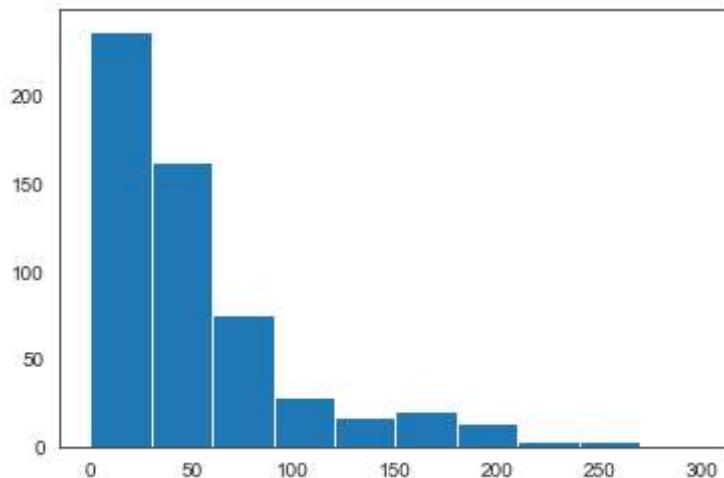


```
In [104]: # <<< chat - 2
```

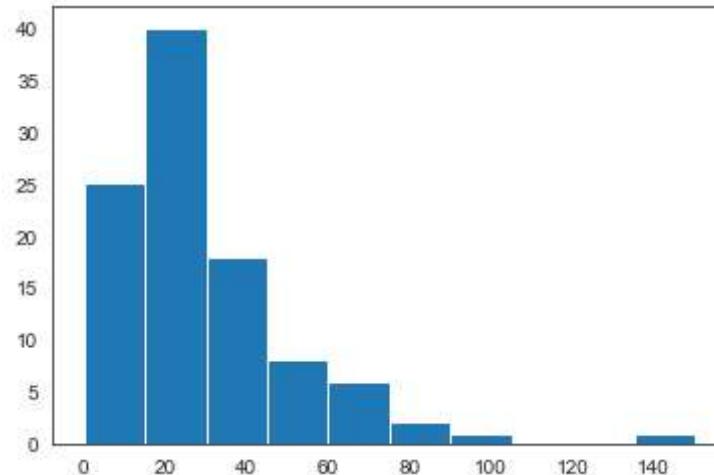
```
# Creating stacked histograms & this is bit tough to understand
```

```
In [105]: #h1 = plt.hist(movies.BudgetMillions)
```

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



```
In [106]: plt.hist(movies[movies.Genre == 'Drama'].BudgetMillions)
plt.show()
```



```
In [107]: movies.head()
```

Out[107]:

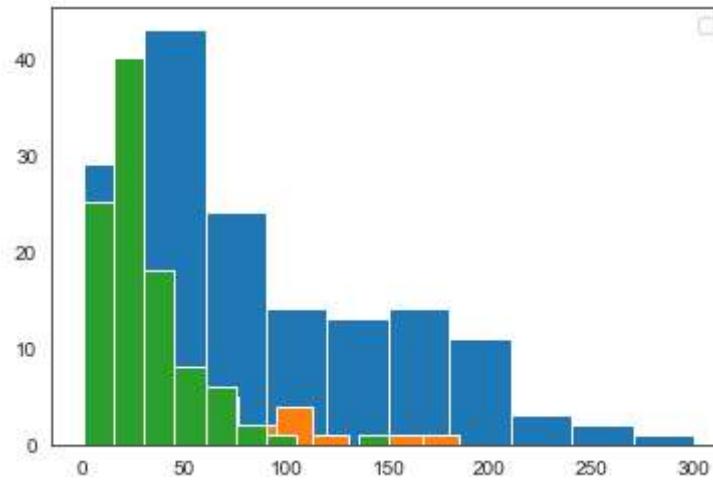
	Film	Genre	CriticRating	AudienceRating	BudgetMillions	Year
0	(500) Days of Summer	Comedy	87	81	8	2009
1	10,000 B.C.	Adventure	9	44	105	2008
2	12 Rounds	Action	30	52	20	2009
3	127 Hours	Adventure	93	84	18	2010
4	17 Again	Comedy	55	70	20	2009

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

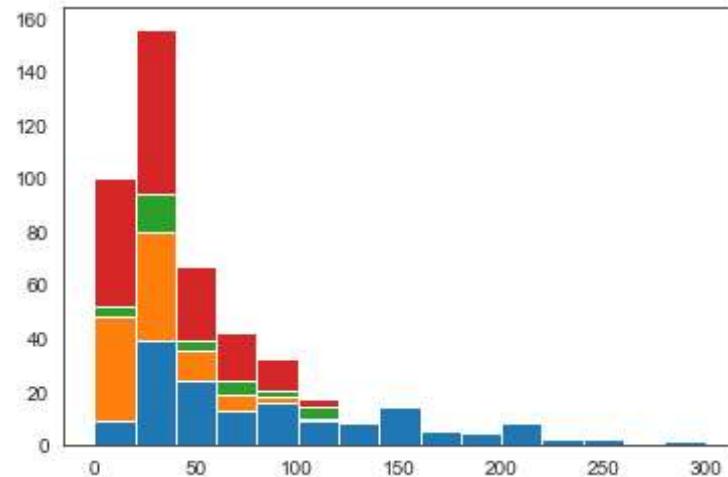
In [109]: # Below plots are stacked histogram because overlaped

```
plt.hist(movies[movies.Genre == 'Action'].BudgetMillions, bins = 10)
plt.hist(movies[movies.Genre == 'Thriller'].BudgetMillions, bins = 10)
plt.hist(movies[movies.Genre == 'Drama'].BudgetMillions, bins = 10)
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 [111]: 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 = 15, stacked = True)\nplt.show()

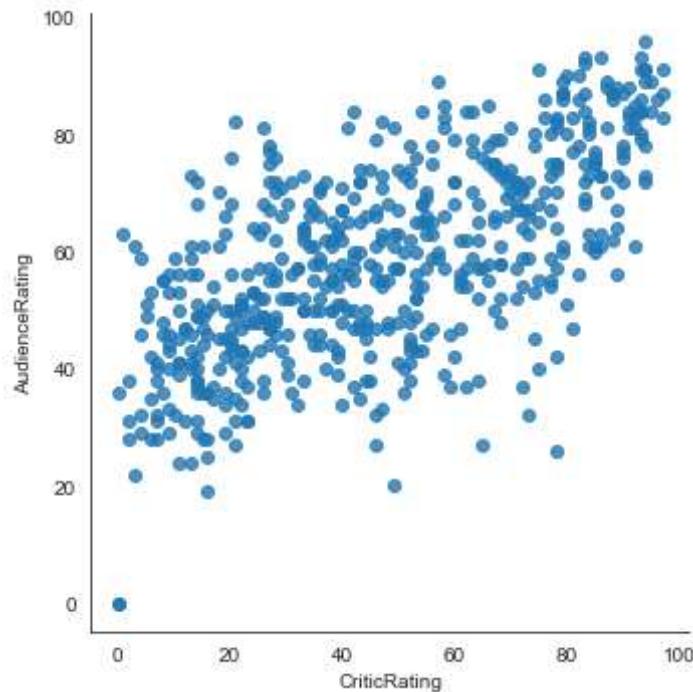


```
In [112]: # if you have 100 categories you cannot copy & paste all the things
```

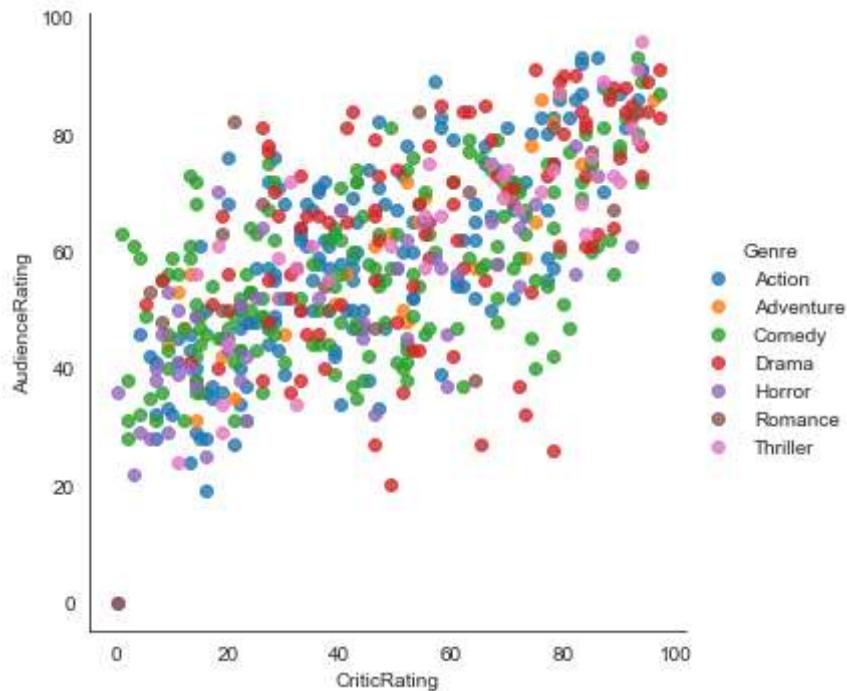
```
for gen in movies.Genre.cat.categories:  
    print(gen)
```

Action
Adventure
Comedy
Drama
Horror
Romance
Thriller

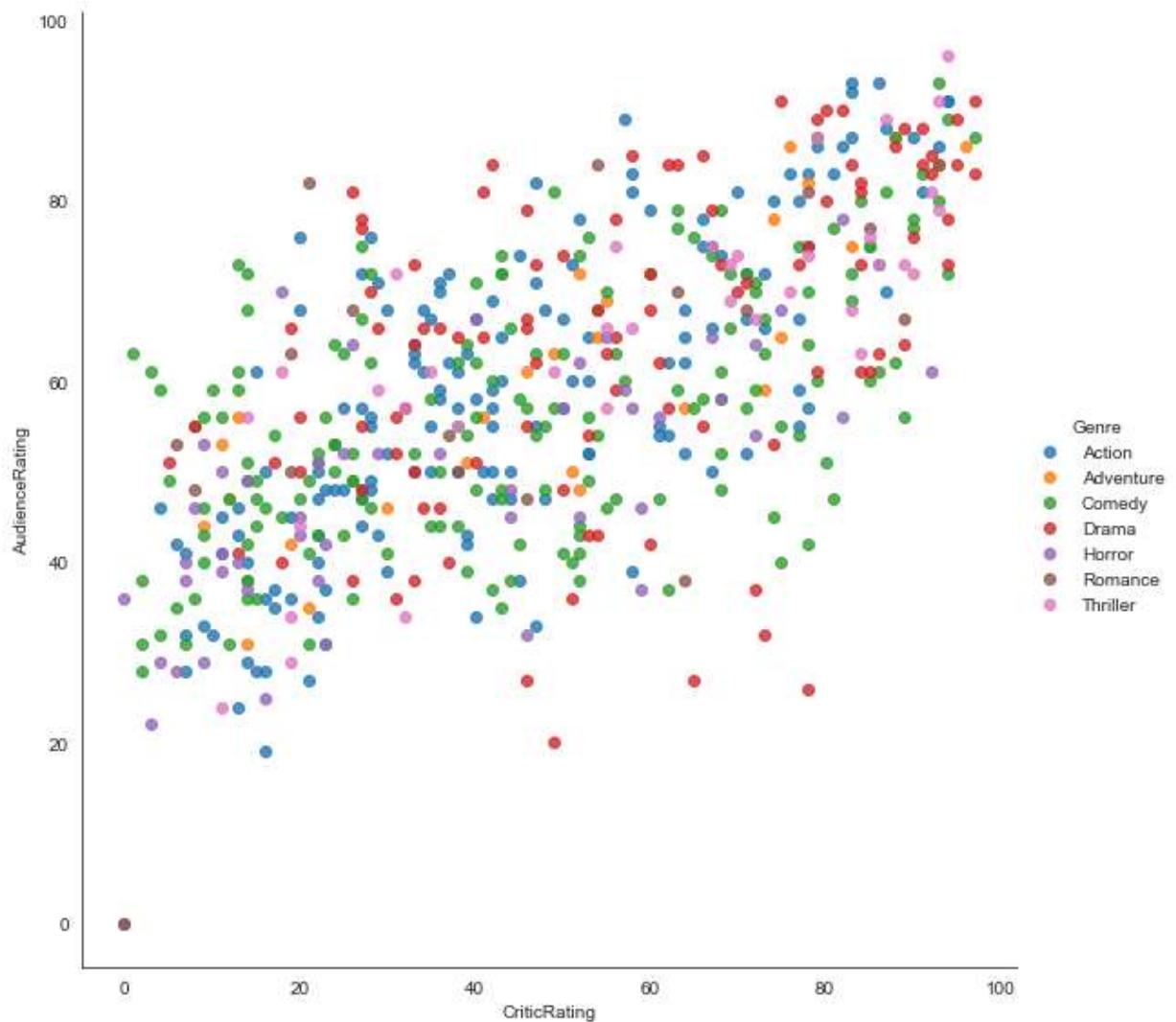
```
In [113]: vis1 = sns.lmplot(data=movies, x='CriticRating', y='AudienceRating',  
                      fit_reg=False)
```



```
In [114]: vis1 = sns.lmplot(data=movies, x='CriticRating', y='AudienceRating',\n                      fit_reg=False, hue = 'Genre')
```



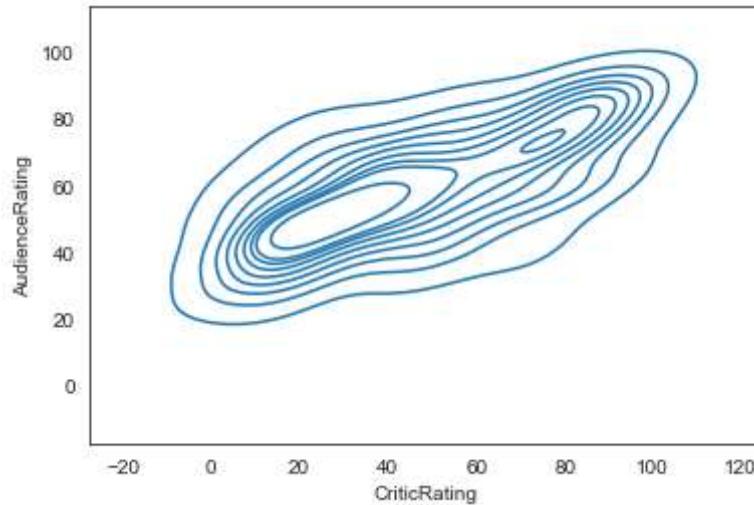
```
In [115]: vis1 = sns.lmplot(data=movies, x='CriticRating', y='AudienceRating',\n                      fit_reg=False, hue = 'Genre', size = 8, aspect=1)
```



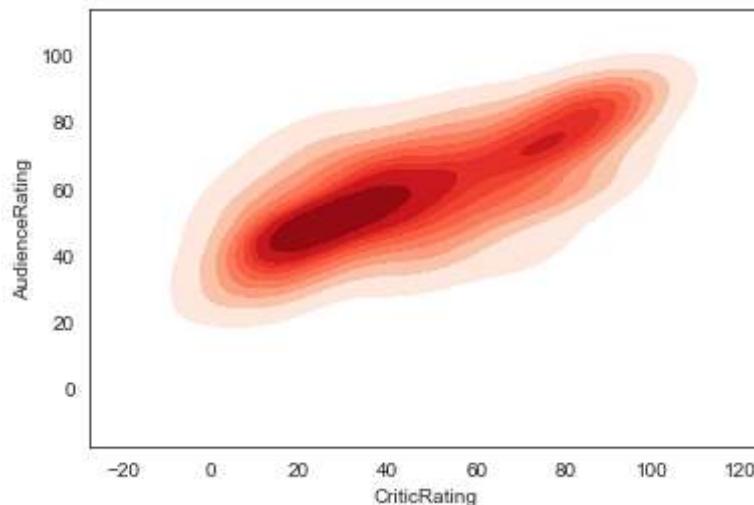
```
In [116]: # Kernel Density Estimate plot ( KDE PLOT)\n# how can i visualize audience rating & critics rating . using scatterplot
```

```
In [117]: k1 = sns.kdeplot(movies.CriticRating,movies.AudienceRating)
```

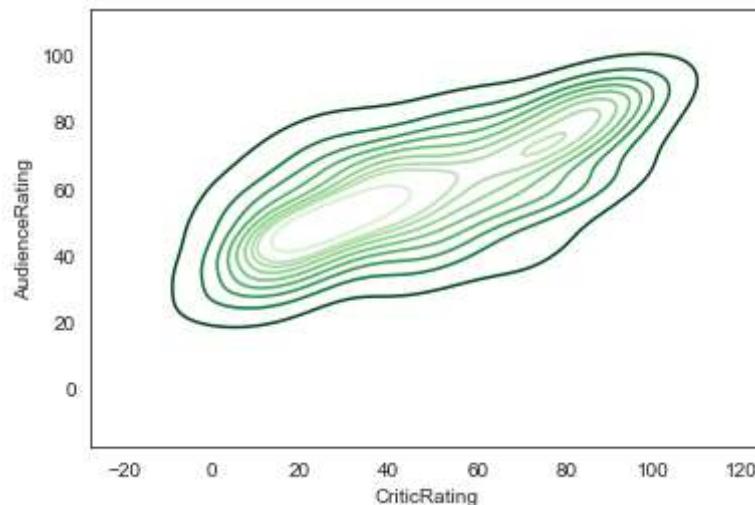
where do u find more density and how density is distributed across the the
center point is kernel this is called KDE & instead of dots it visualize like t
we can able to clearly see the spread at the audience ratings



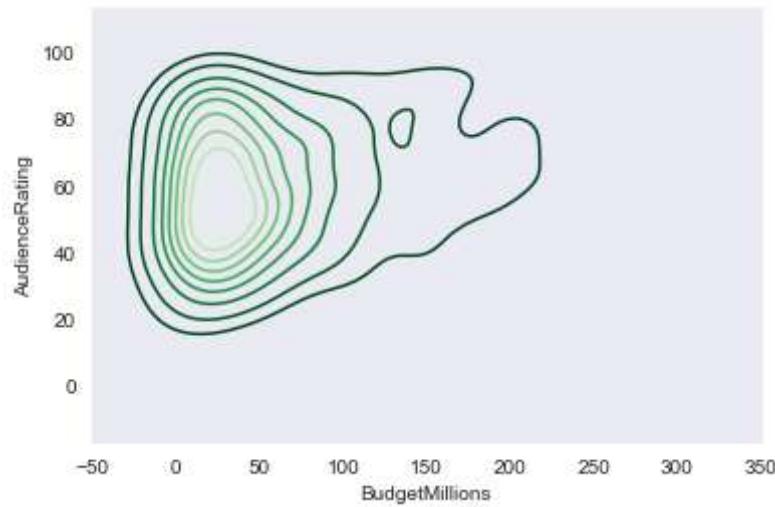
```
In [118]: k1 = sns.kdeplot(movies.CriticRating,movies.AudienceRating,shade = True,shade_low
```



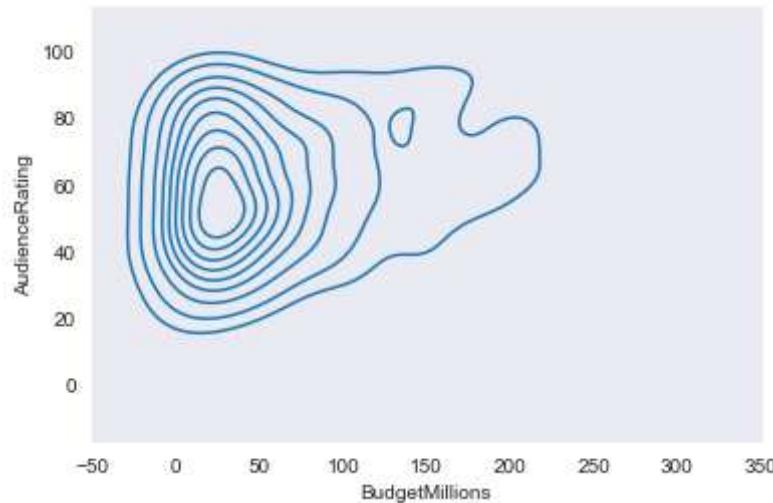
```
In [119]: k2 = sns.kdeplot(movies.CriticRating,movies.AudienceRating,shade_lowest=False,cmap="viridis")
```



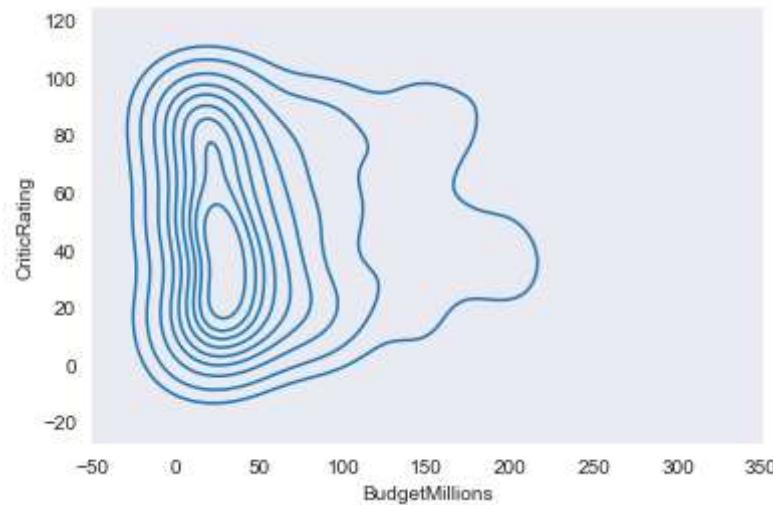
```
In [120]: sns.set_style('dark')
k1 = sns.kdeplot(movies.BudgetMillions,movies.AudienceRating,shade_lowest=False,cmap="viridis")
```



```
In [121]: sns.set_style('dark')
k1 = sns.kdeplot(movies.BudgetMillions,movies.AudienceRating)
```

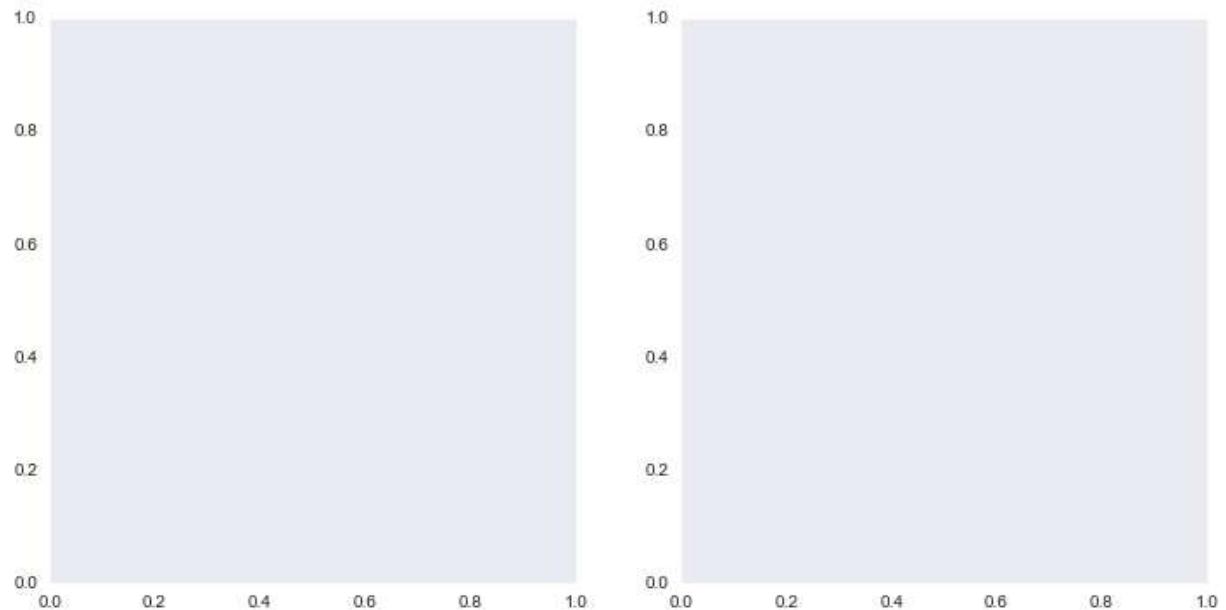


```
In [122]: k2 = sns.kdeplot(movies.BudgetMillions,movies.CriticRating)
```



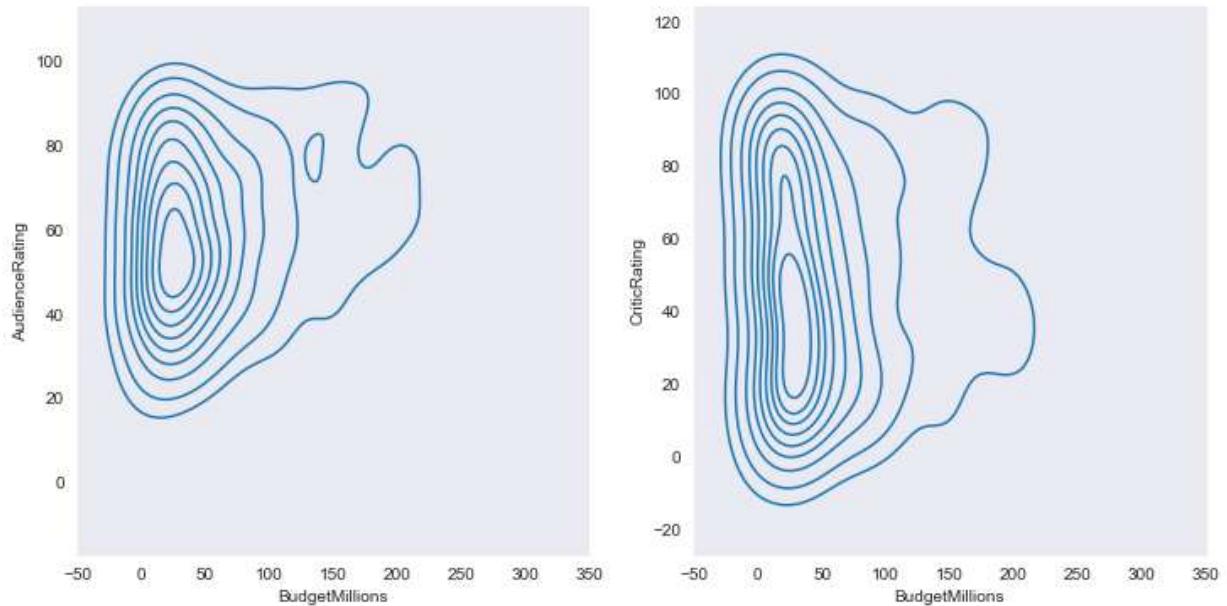
In [123]: `#subplots`

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



```
In [124]: f, axes = plt.subplots(1,2, figsize =(12,6))

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

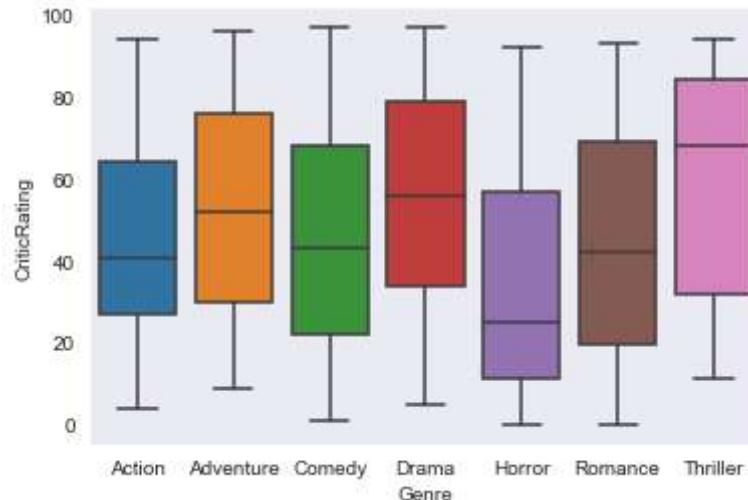


```
In [125]: axes
```

```
Out[125]: array([<AxesSubplot:xlabel='BudgetMillions', ylabel='AudienceRating'>,
   <AxesSubplot:xlabel='BudgetMillions', ylabel='CriticRating'>],  
dtype=object)
```

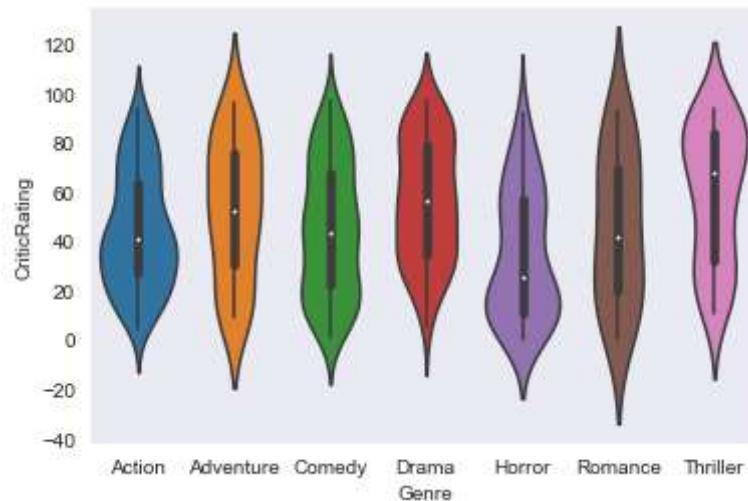
In [126]: #Box plots -

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

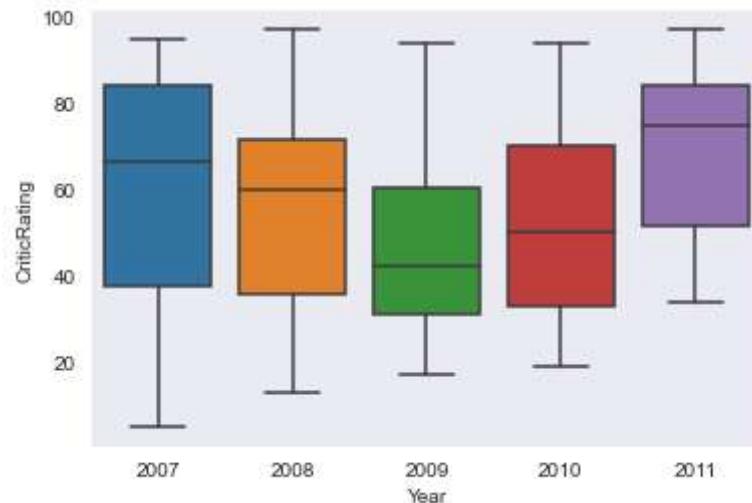


In [127]: #violin plot

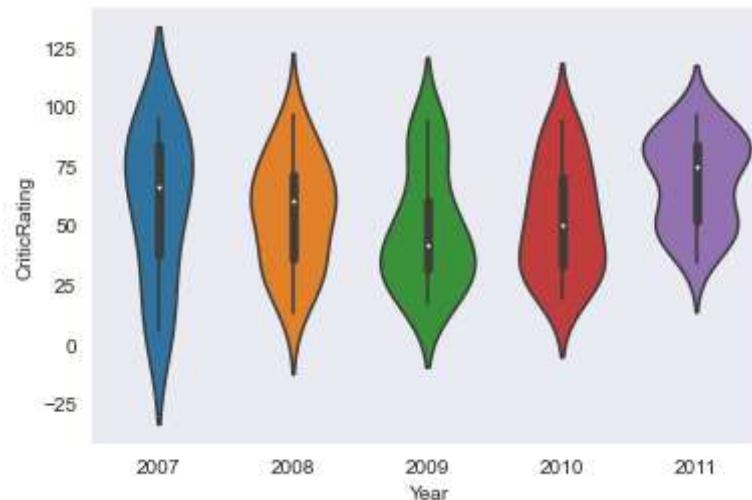
```
z = sns.violinplot(data=movies, x='Genre', y = 'CriticRating')
```



```
In [128]: w1 = sns.boxplot(data=movies[movies.Genre == 'Drama'], x='Year', y = 'CriticRating')
```

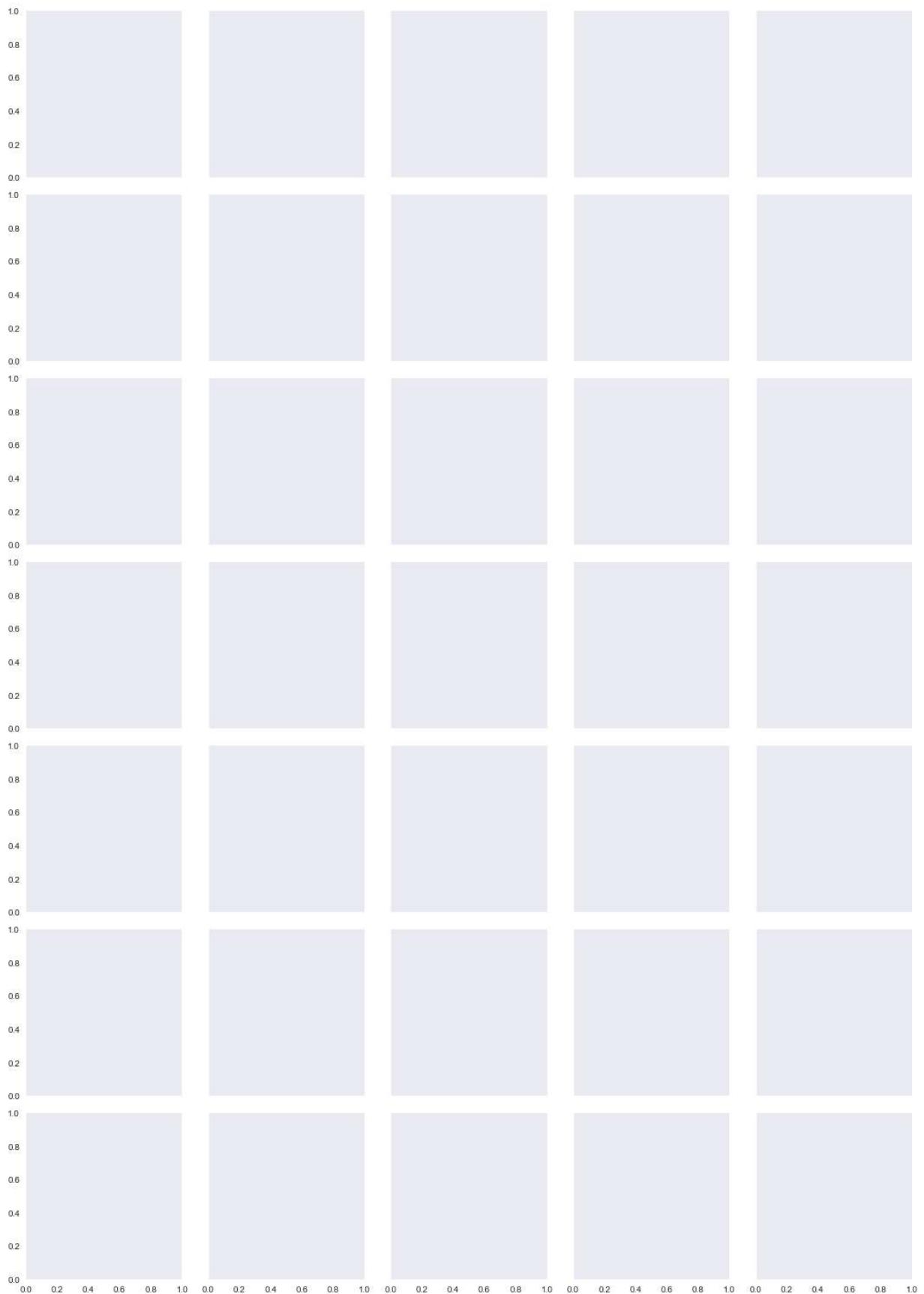


```
In [129]: z = sns.violinplot(data=movies[movies.Genre == 'Drama'], x='Year', y = 'CriticRating')
```



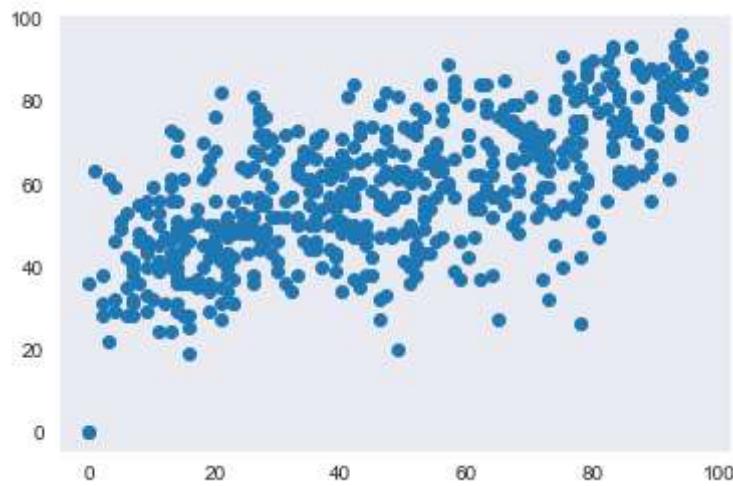
```
In [130]: # Creating a Facet grid
```

```
In [131]: g = sns.FacetGrid (movies, row = 'Genre', col = 'Year', hue = 'Genre') #kind of su
```



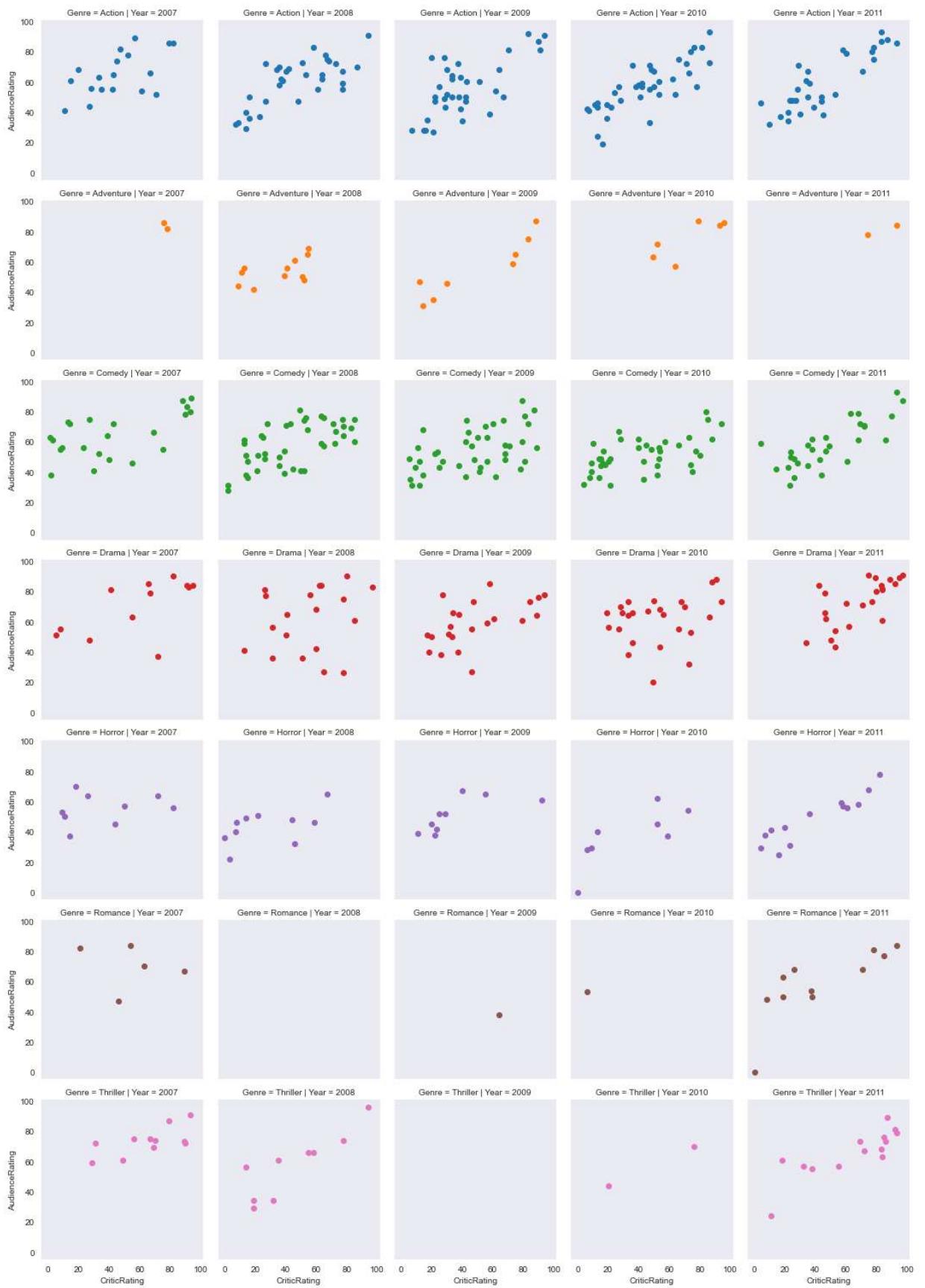
```
In [132]: plt.scatter(movies.CriticRating,movies.AudienceRating)
```

```
Out[132]: <matplotlib.collections.PathCollection at 0x16141e2f8e0>
```



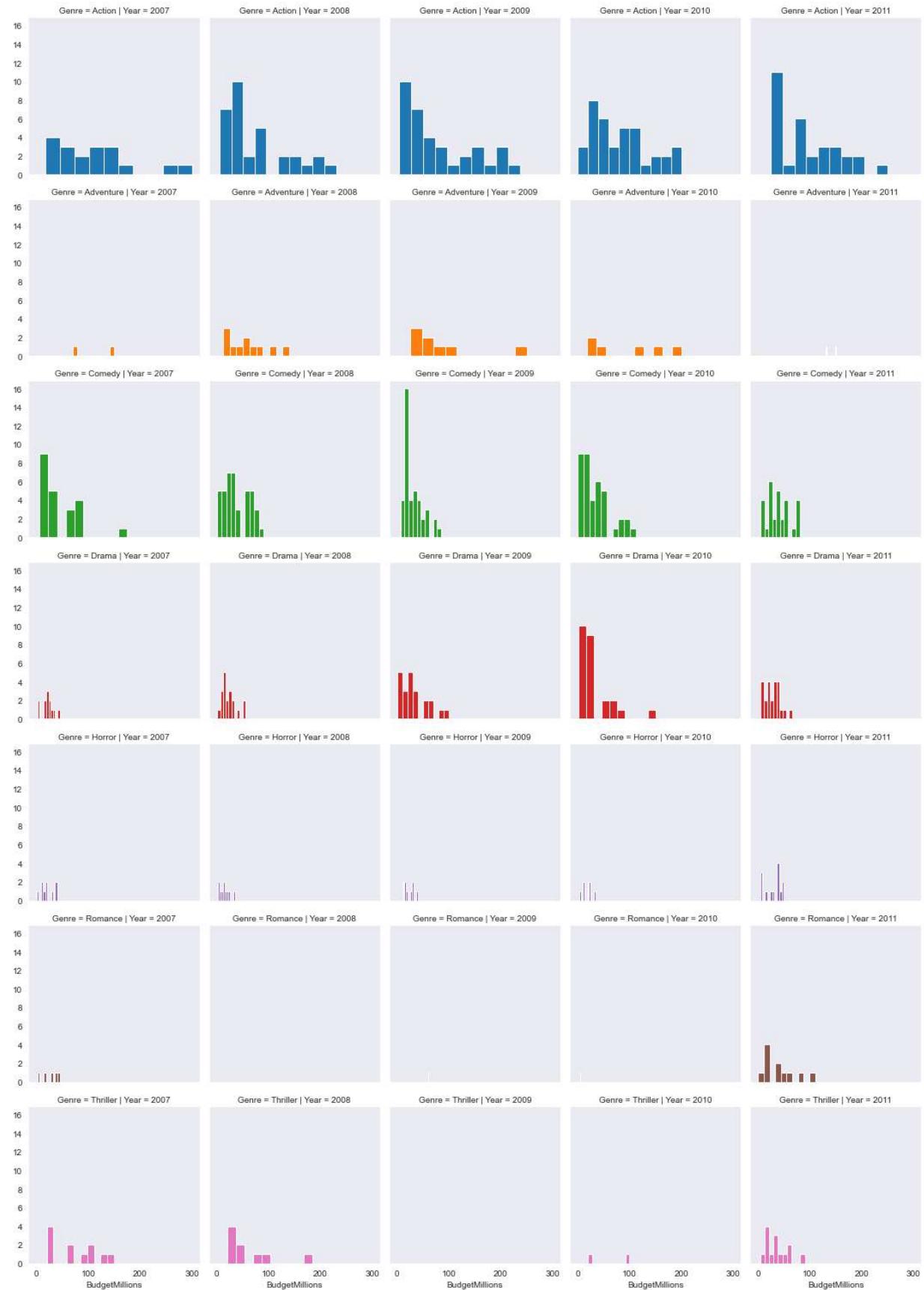
In [133]:

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



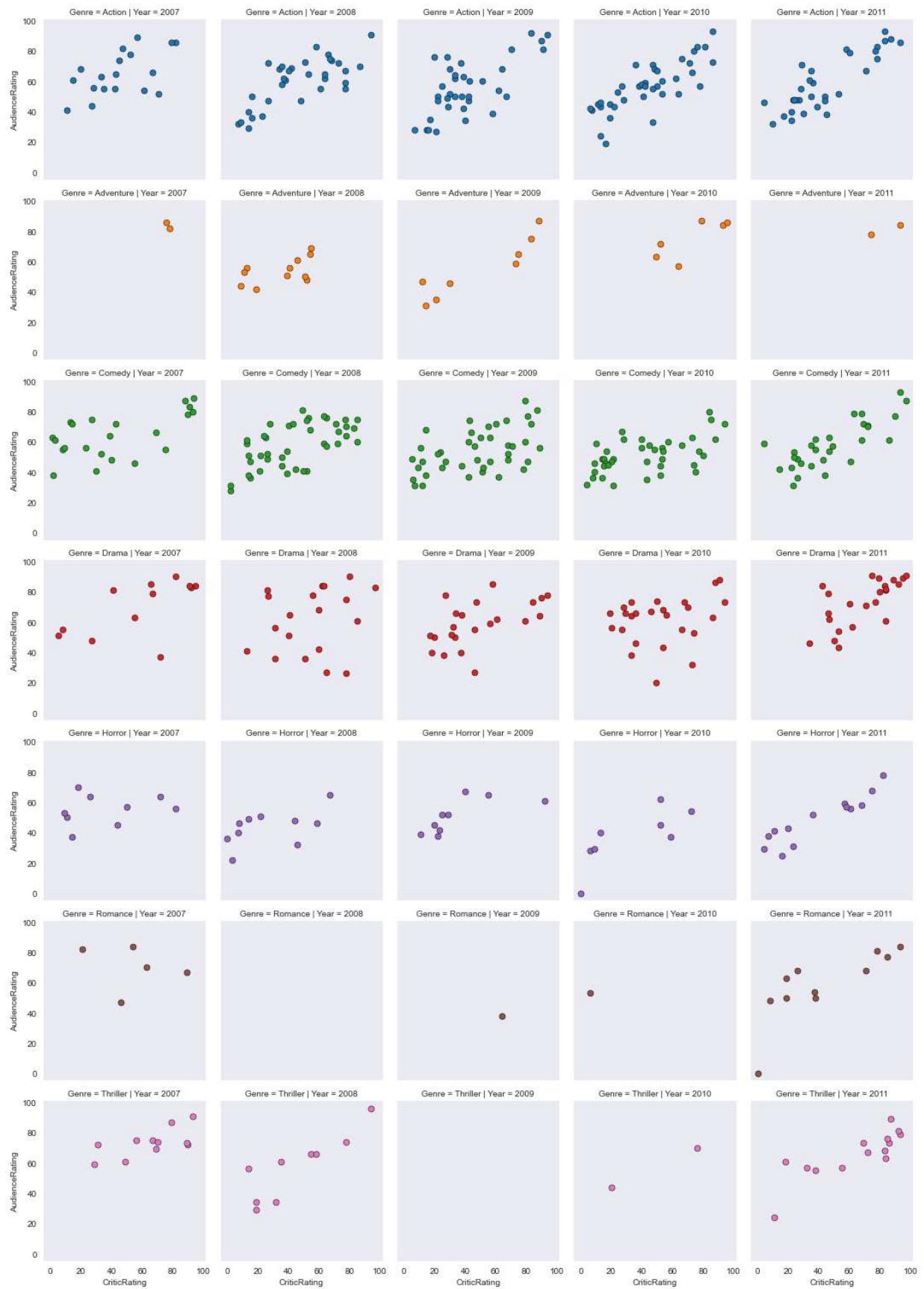
In [134]: # you can populated any type of chat.

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



In [135]:

```
#  
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 [136]: # 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))

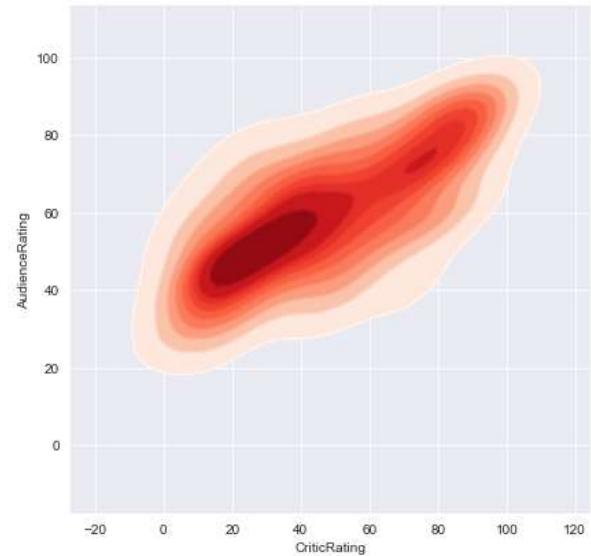
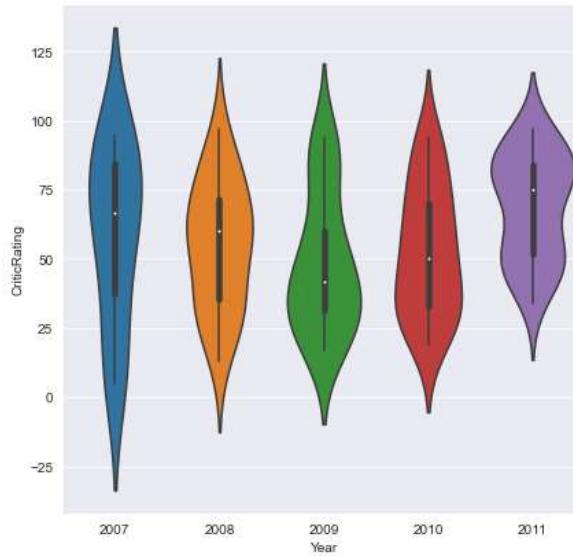
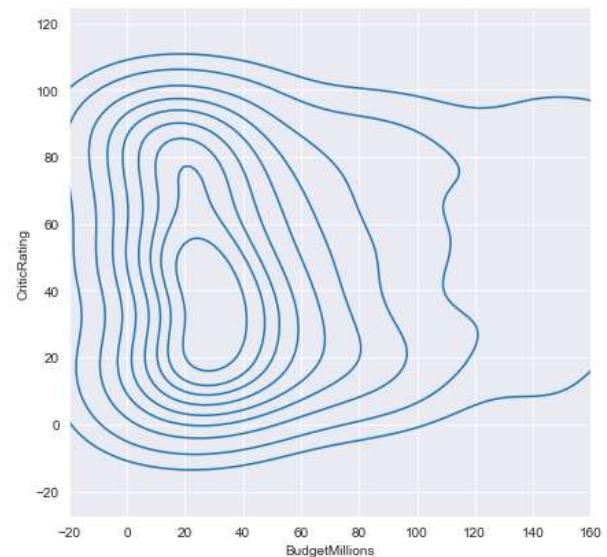
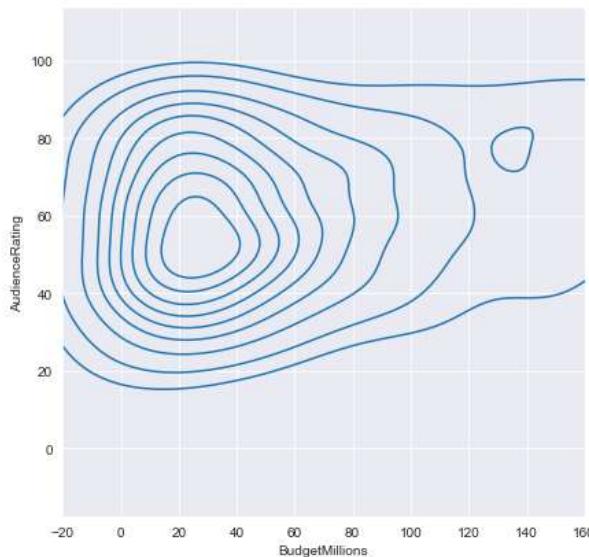
k1 = sns.kdeplot(movies.BudgetMillions,movies.AudienceRating,ax=axes[0,0])
k2 = sns.kdeplot(movies.BudgetMillions,movies.CriticRating,ax = axes[0,1])

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(movies.CriticRating,movies.AudienceRating,shade = True,shade_lowest=False)
k4b = sns.kdeplot(movies.CriticRating, movies.AudienceRating,cmap='Reds',ax = axes[1,0])

plt.show()
```



In [144]: # How can you style your dashboard using different color map

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

sns.set_style('dark',{'axes.facecolor':'black'})
f, axes = plt.subplots(2,2, figsize = (15,15))

#plot [0,0]
k1 = sns.kdeplot(movies.BudgetMillions,movies.AudienceRating, \
                  shade = True, shade_lowest=True,cmp = 'inferno', \
                  ax = axes[0,0])
k1b = sns.kdeplot(movies.BudgetMillions, movies.AudienceRating, \
                   cmap = 'cool',ax = axes[0,0])

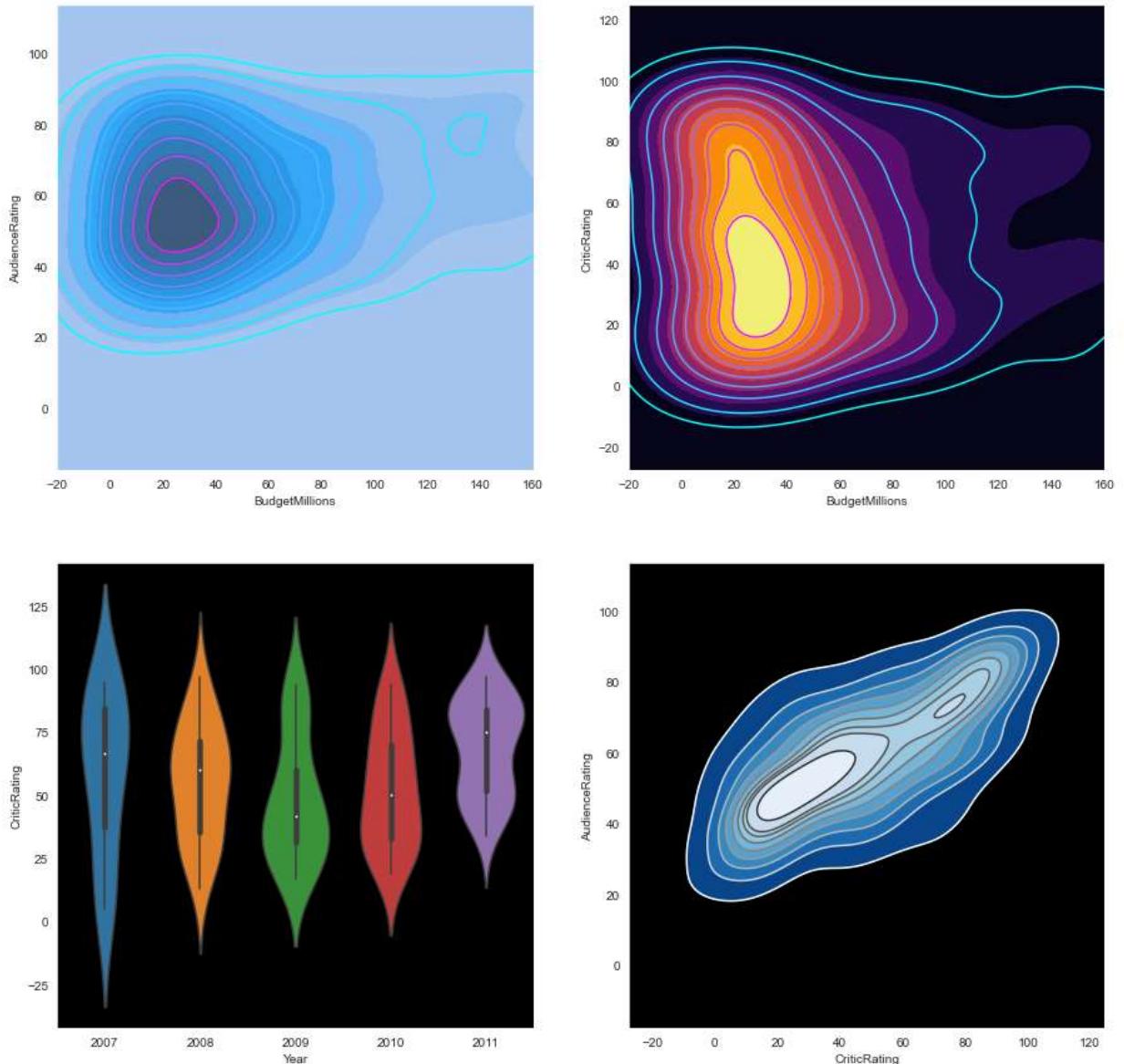
#plot [0,1]
k2 = sns.kdeplot(movies.BudgetMillions,movies.CriticRating, \
                  shade=True, shade_lowest=True, cmap='inferno', \
                  ax = axes[0,1])
k2b = sns.kdeplot(movies.BudgetMillions,movies.CriticRating, \
                   cmap = 'cool', ax = axes[0,1])

#plot[1,0]
z = sns.violinplot(data=movies[movies.Genre=='Drama'], \
                     x='Year', y = 'CriticRating', ax=axes[1,0])

#plot[1,1]
k4 = sns.kdeplot(movies.CriticRating,movies.AudienceRating, \
                  shade = True,shade_lowest=False,cmap='Blues_r', \
                  ax=axes[1,1])
k4b = sns.kdeplot(movies.CriticRating, movies.AudienceRating, \
                   cmap='gist_gray_r',ax = axes[1,1])

k1.set(xlim=(-20,160))
k2.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

In [145]: # eda is completed

In []: