

In [134]:

```
1 import pandas as pd
2 import os
```

In [135]:

```
1 os.getcwd()
```

Out[135]:

```
'/Users/myyntiimac'
```

In [136]:

```
1 df=pd.read_csv("/Users/myyntiimac/Desktop/P4-Movie-Ratings.csv")
2 df.head()
```

Out[136]:

	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 [145]:

```
1 len(df)
```

Out[145]:

```
559
```

In [ ]:

```
1
```

In [146]:

```
1 df.shape
```

Out[146]:

```
(559, 6)
```

In [147]:

```
1 df.columns
```

Out[147]:

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

In [148]:

```
1 df.columns=['Film', 'Genre', 'Criticsratings', 'AudienceRatings',
2           'Budgetmillion', 'Year']
```

In [149]:

```
1 df.head()
```

Out[149]:

	Film	Genre	Criticsratings	AudienceRatings	Budgetmillion	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 [150]:

```
1 df.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   Criticsratings        559 non-null    int64
3   AudienceRatings       559 non-null    int64
4   Budgetmillion         559 non-null    int64
5   Year                  559 non-null    int64
dtypes: int64(4), object(2)
memory usage: 26.3+ KB
```

In [151]:

```
1 df.isnull().any()
```

Out[151]:

```
Film           False
Genre          False
Criticsratings False
AudienceRatings False
Budgetmillion  False
Year           False
dtype: bool
```

In [152]:

```
1 df.describe().transpose()
```

Out[152]:

	count	mean	std	min	25%	50%	75%	max
<b>Criticsratings</b>	559.0	47.309481	26.413091	0.0	25.0	46.0	70.0	97.0
<b>AudienceRatings</b>	559.0	58.744186	16.826887	0.0	47.0	58.0	72.0	96.0
<b>Budgetmillion</b>	559.0	50.236136	48.731817	0.0	20.0	35.0	65.0	300.0
<b>Year</b>	559.0	2009.152057	1.362632	2007.0	2008.0	2009.0	2010.0	2011.0

In [ ]:

```
1 #we dont want year treat as number, we can make it catagorical
2 #in addition , film and genre showing as object , thats also need to converta as
3
```

In [153]:

```
1 df['Film'] = df['Film'].astype('category')
2
```

In [154]:

```
1 df.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   Criticsratings        559 non-null   int64
3   AudienceRatings       559 non-null   int64
4   Budgetmillion         559 non-null   int64
5   Year                  559 non-null   int64
dtypes: category(1), int64(4), object(1)
memory usage: 43.6+ KB
```

In [155]:

```
1 df['Genre'] = df['Genre'].astype('category')
```

In [156]:

```
1 df.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   Criticsratings        559 non-null   int64
3   AudienceRatings      559 non-null   int64
4   Budgetmillion         559 non-null   int64
5   Year                  559 non-null   int64
dtypes: category(2), int64(4)
memory usage: 40.1 KB
```

In [21]:

```
1 df.Year=df.Year.astype("category")
```

In [157]:

```
1 df.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   Criticsratings        559 non-null   int64
3   AudienceRatings      559 non-null   int64
4   Budgetmillion         559 non-null   int64
5   Year                  559 non-null   int64
dtypes: category(2), int64(4)
memory usage: 40.1 KB
```

In [158]:

```
1 #to asses unique values in categories
2 df.Genre.cat.categories
```

Out[158]:

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

In [159]:

```
1 df.Genre.unique()
```

Out[159]:

```
['Comedy', 'Adventure', 'Action', 'Horror', 'Drama', 'Romance', 'Thriller']
Categories (7, object): ['Action', 'Adventure', 'Comedy', 'Drama', 'Horror', 'Romance', 'Thriller']
```

In [160]:

```
1 df.describe()
```

Out[160]:

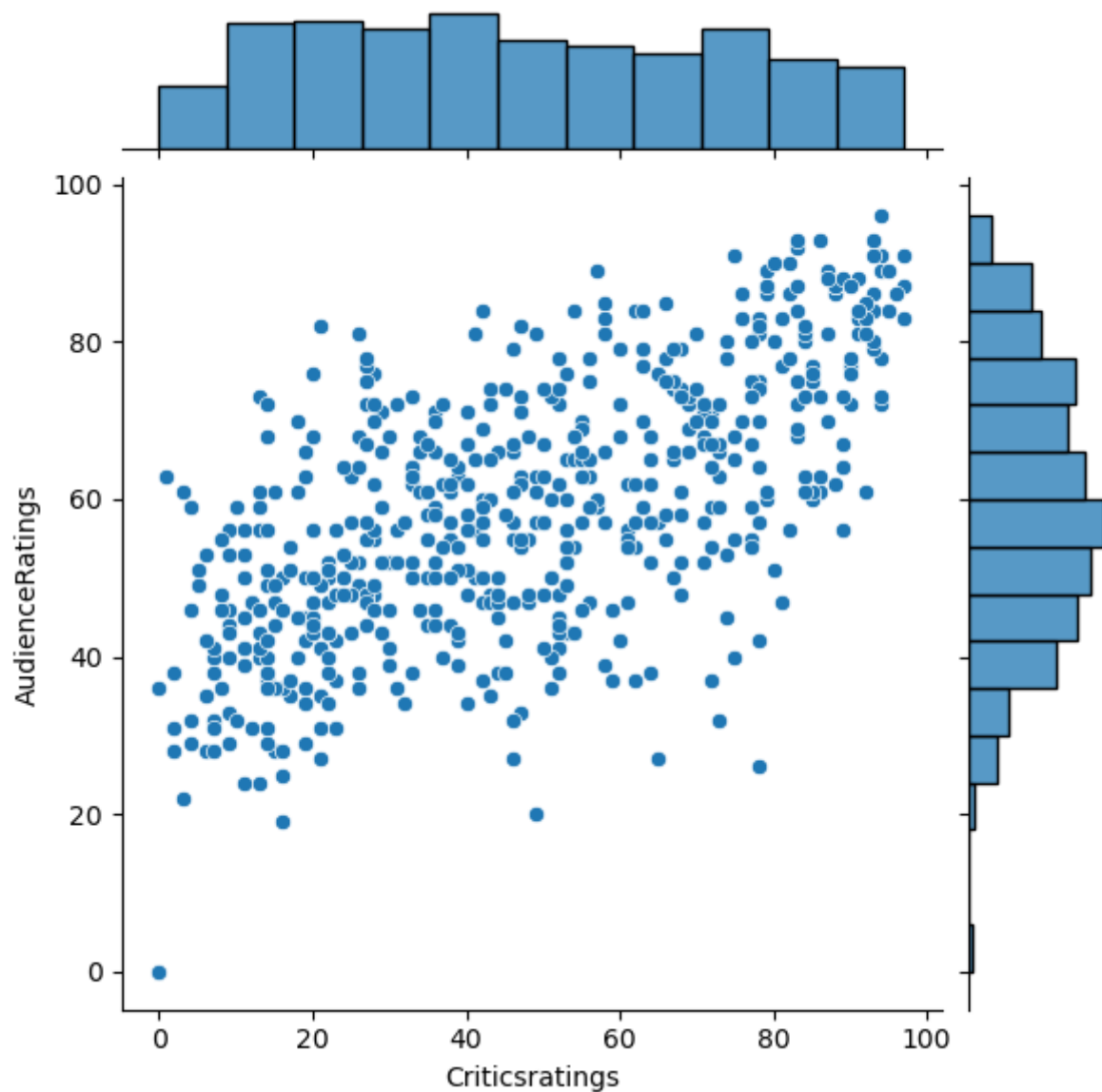
	Criticsratings	AudienceRatings	Budgetmillion	Year
<b>count</b>	559.000000	559.000000	559.000000	559.000000
<b>mean</b>	47.309481	58.744186	50.236136	2009.152057
<b>std</b>	26.413091	16.826887	48.731817	1.362632
<b>min</b>	0.000000	0.000000	0.000000	2007.000000
<b>25%</b>	25.000000	47.000000	20.000000	2008.000000
<b>50%</b>	46.000000	58.000000	35.000000	2009.000000
<b>75%</b>	70.000000	72.000000	65.000000	2010.000000
<b>max</b>	97.000000	96.000000	300.000000	2011.000000

In [106]:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4 import warnings
5 warnings.filterwarnings("ignore")
```

In [33]:

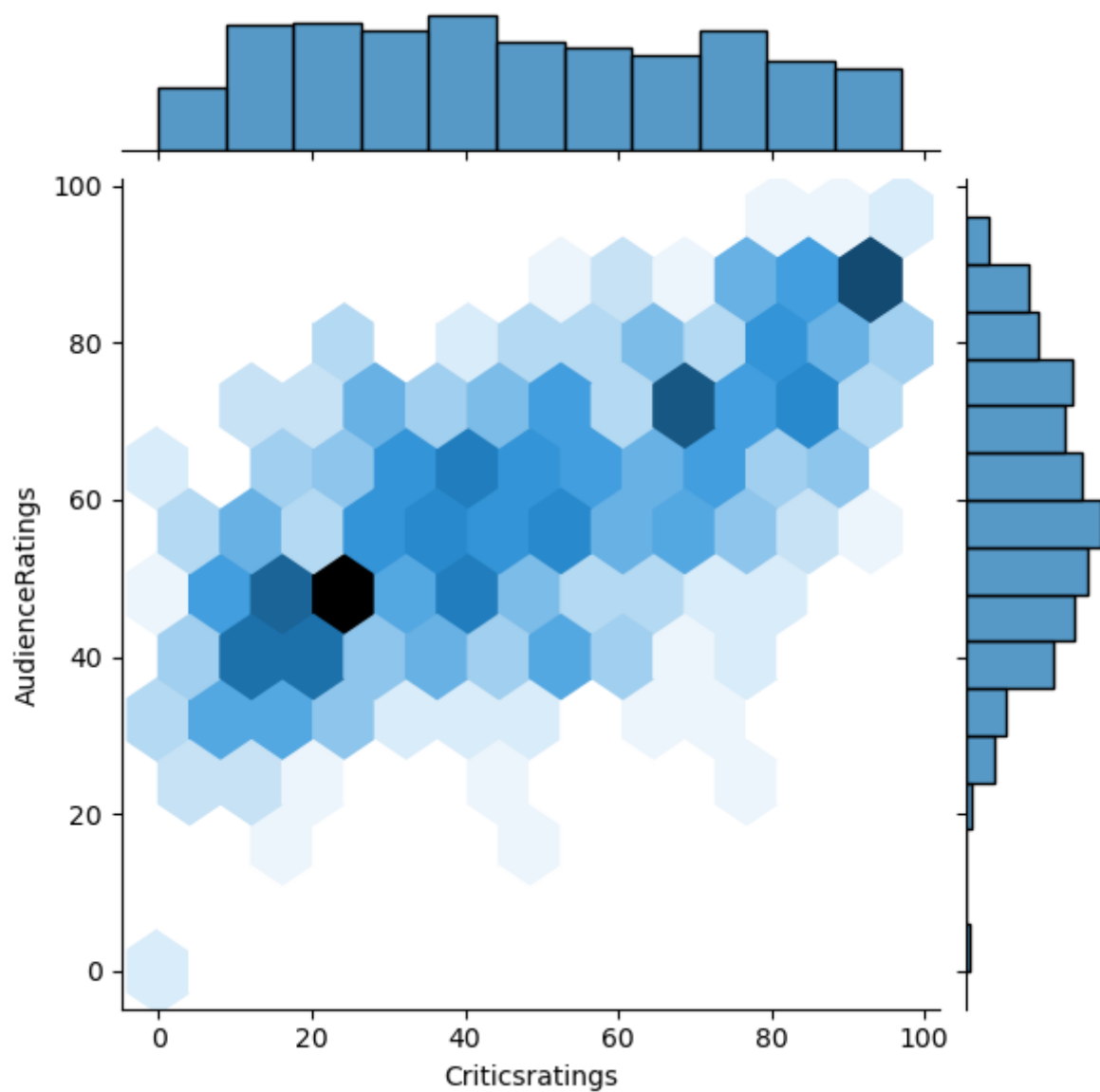
```
1 #creating jointplot
2 a=sns.jointplot(data=df, x="Criticsratings",y="AudienceRatings")
```



Insight: joint plot, you can visualize both univariate and bivariate in scatterplot, you see the relationship between criticsratings and audience ratings and histogram you see the distribution, in distribution we saw audience is normally distributed so it's a good predictor

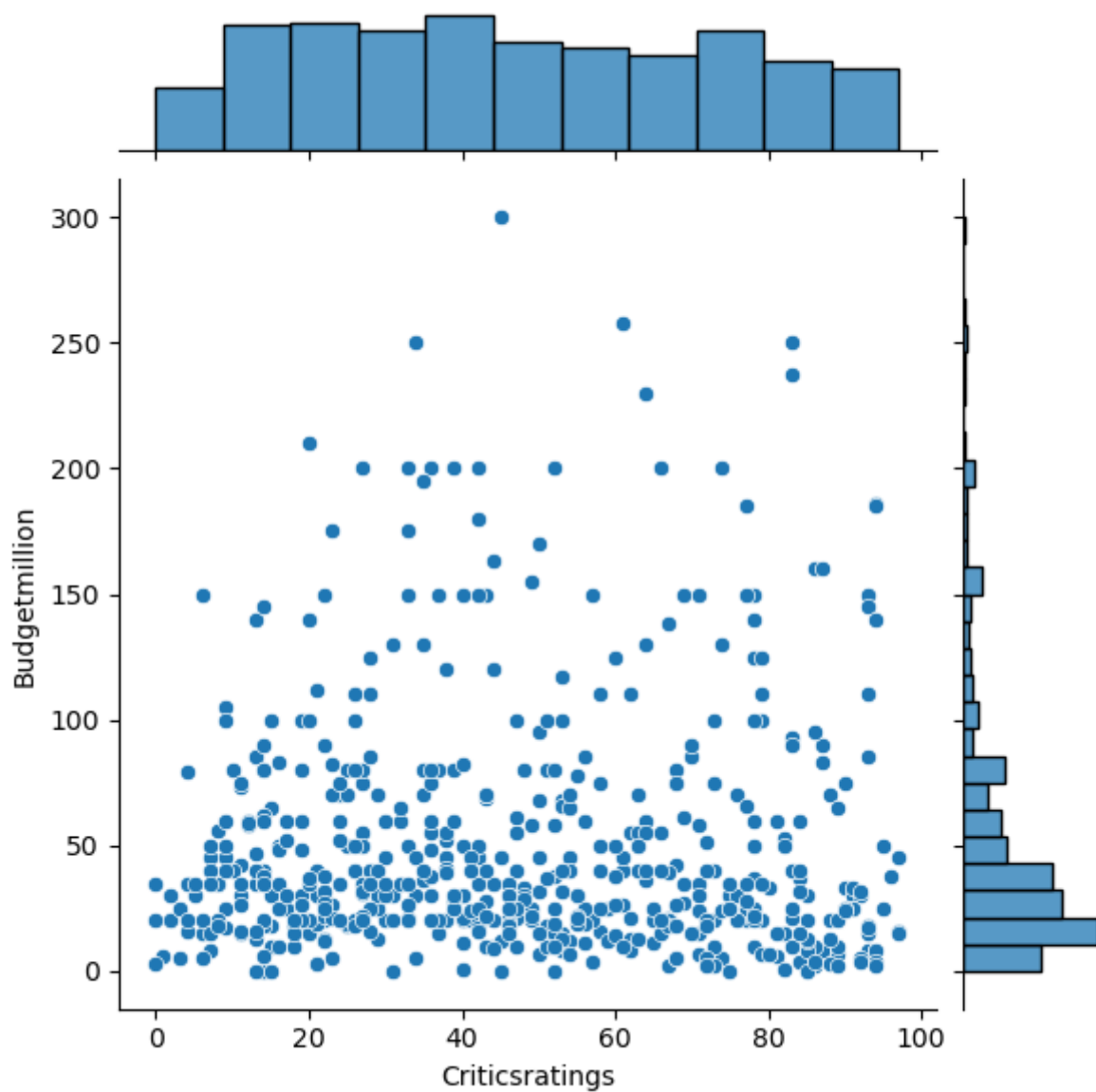
In [34]:

```
1 a=sns.jointplot(data=df, x="Criticsratings",y="AudienceRatings")
```



In [35]:

```
1 a=sns.jointplot(data=df, x="Criticsratings",y="Budgetmillion")
```



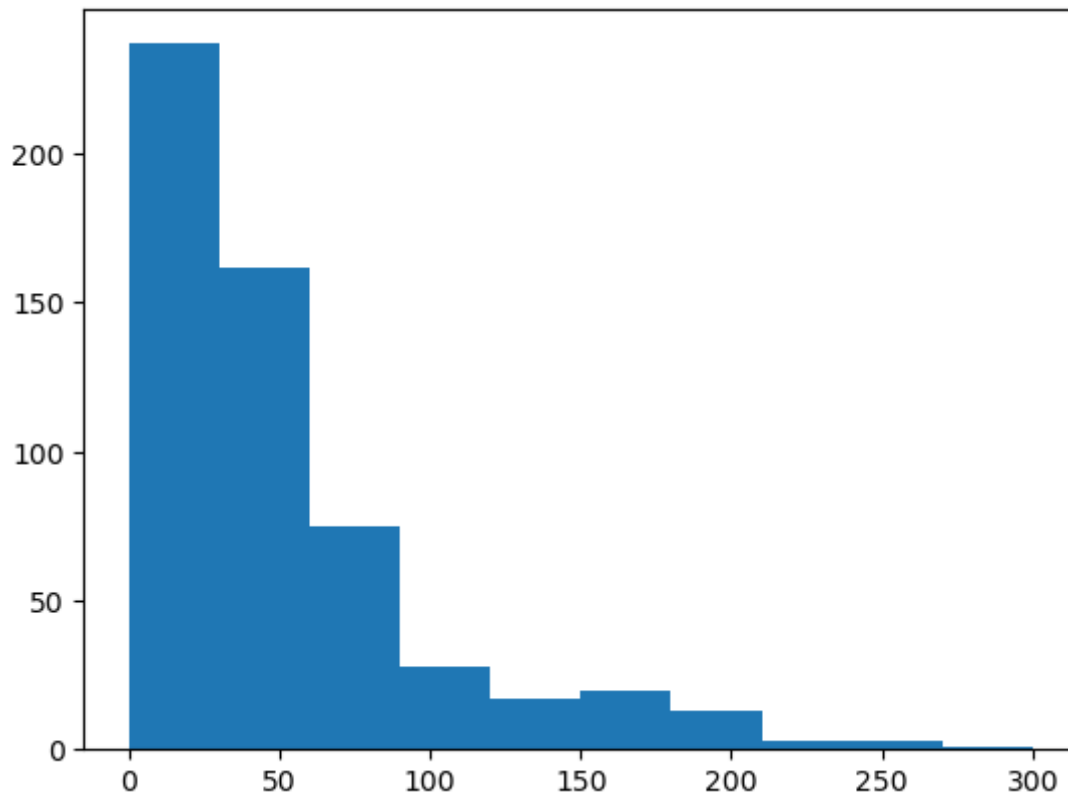


In [38]:

```

1 # creating stack distribution plot
2 plt.hist(df.h=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
3 h=h.map(plt.hist,"Criticsratings"))
4 plt.show()

```



In [39]:

```

1 #create a new data frame contain only Genre= Drama and their budget
2 df[df.Genre=="Drama"].Budgetmillion

```

Out[39]:

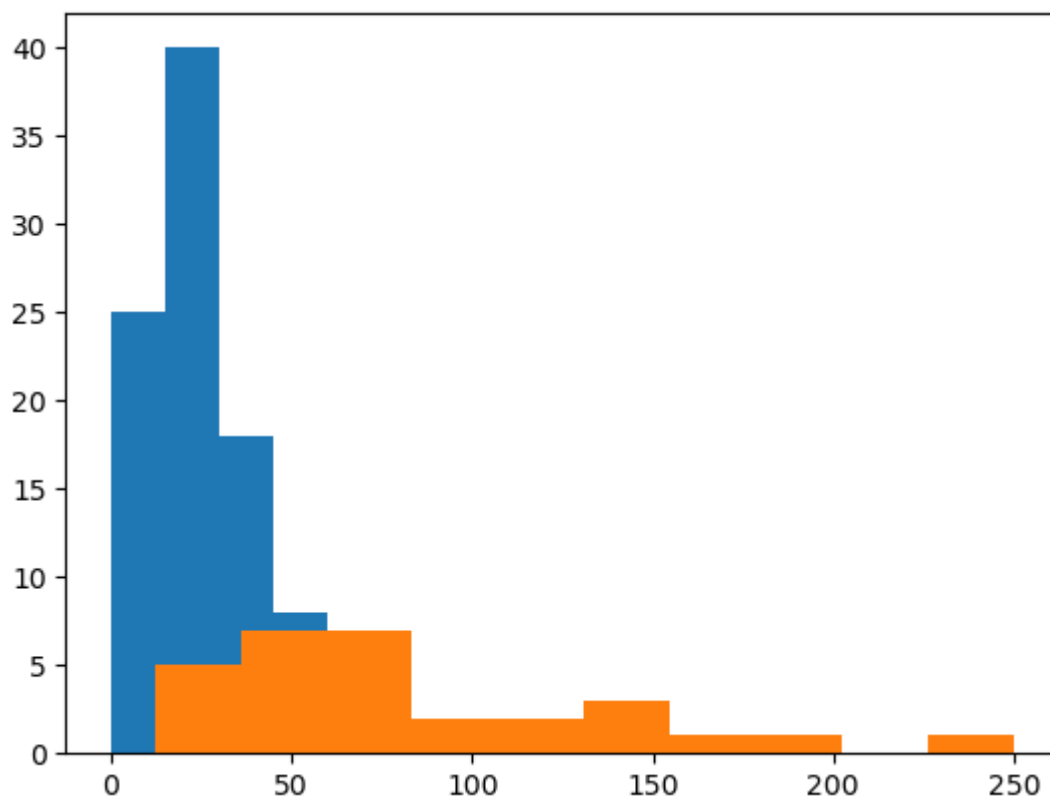
```

10      30
11      20
13       7
18       8
23      20
..
529     66
532     38
534     21
541     15
545       2
Name: Budgetmillion, Length: 101, dtype: int64

```

In [42]:

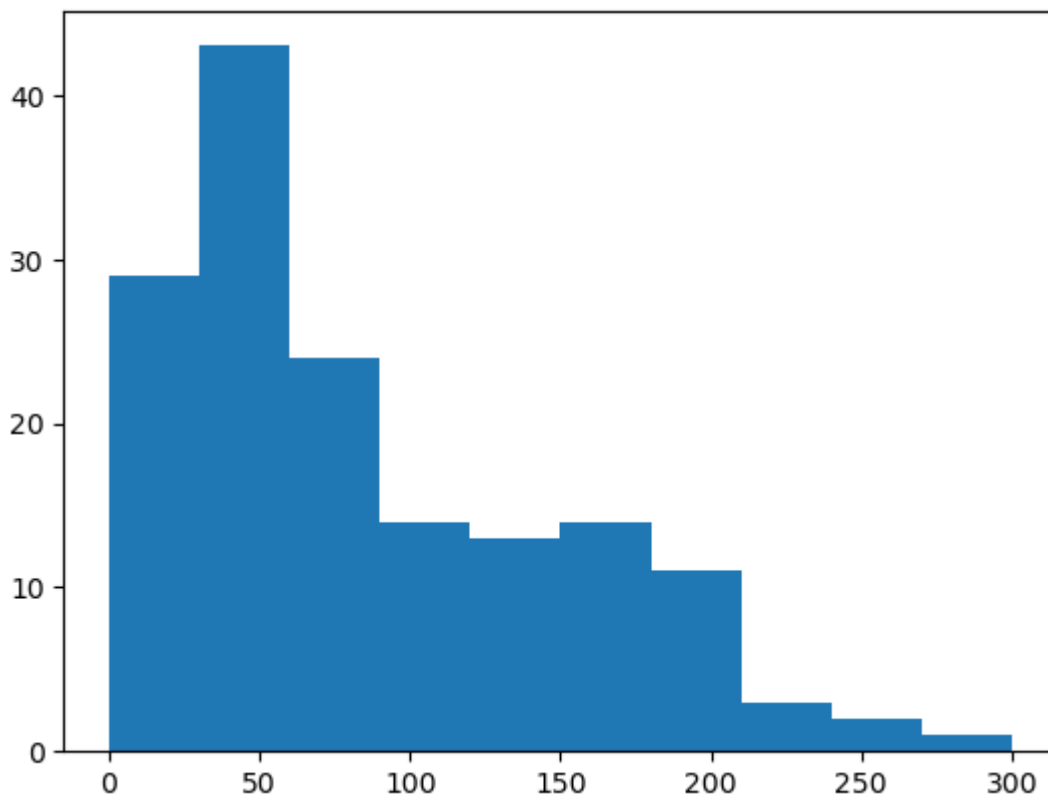
```
1 plt.hist(df[df. Genre=="Drama"].Budgetmillion)
2 plt.hist(df[df. Genre=="Adventure"].Budgetmillion)
3 plt.show()
```



Insight its only for Drama genre and their budget

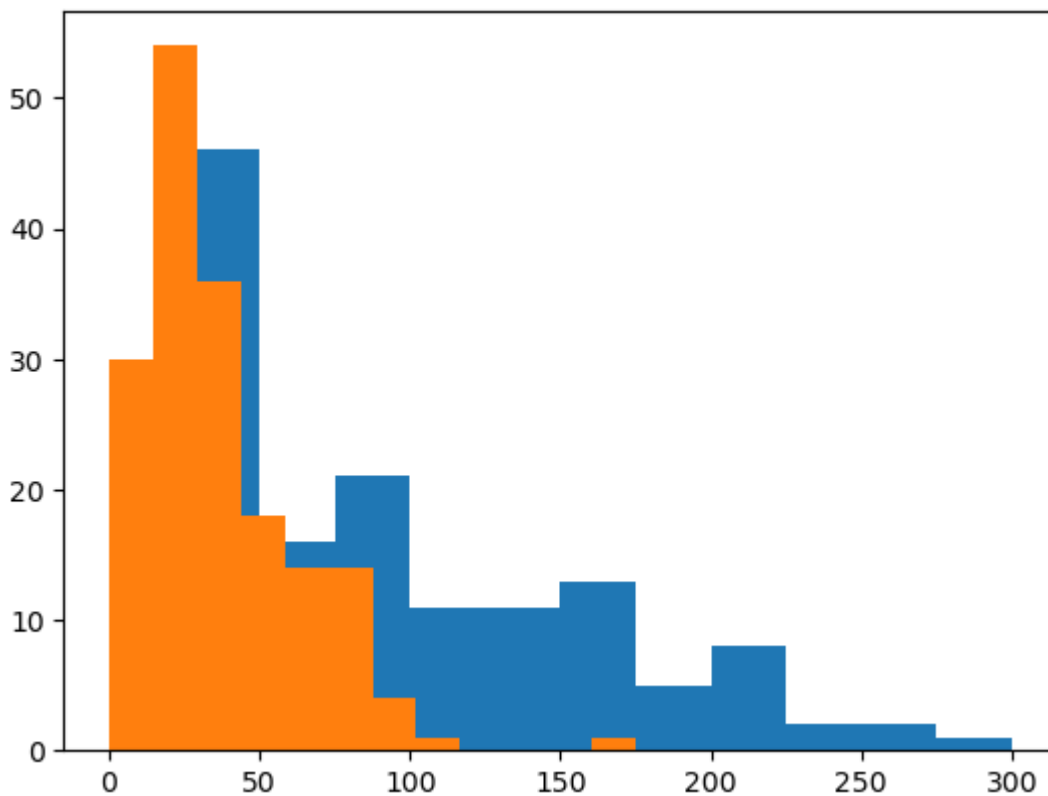
In [41]:

```
1 plt.hist(df[df. Genre=="Action"].Budgetmillion)
2 plt.show()
```



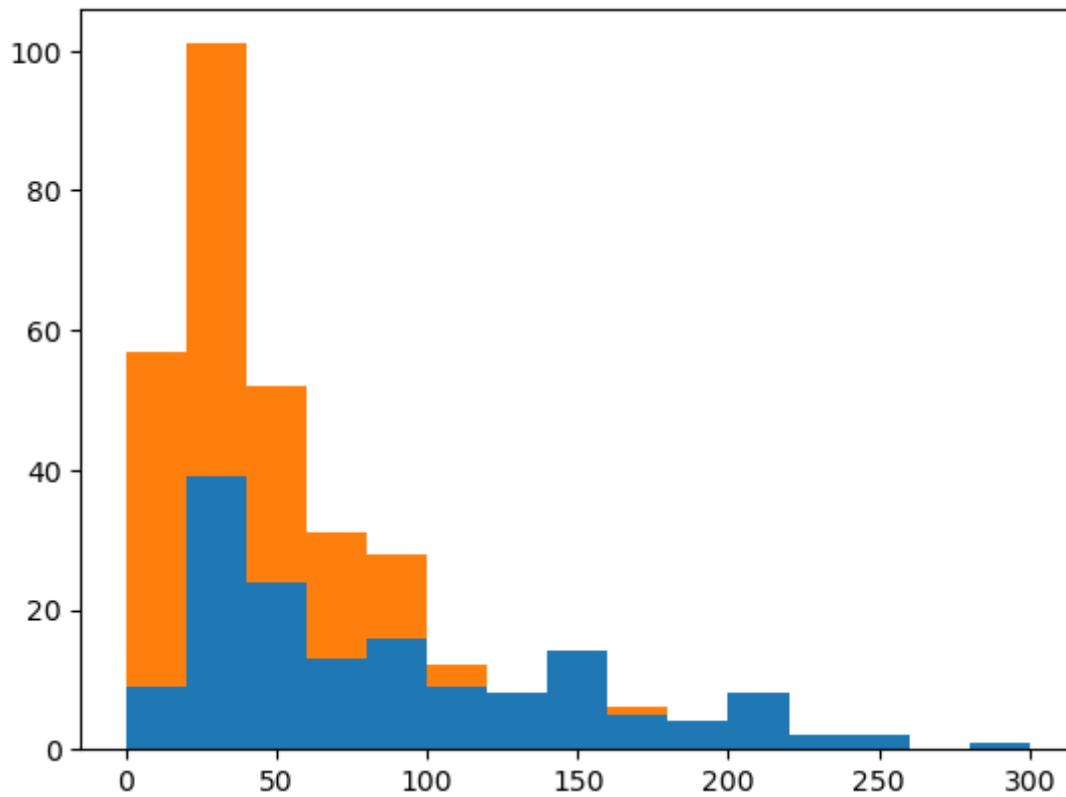
In [44]:

```
1 plt.hist(df[df. Genre=="Action"].Budgetmillion,bins=12)
2 plt.hist(df[df. Genre=="Comedy"].Budgetmillion,bins=12)
3 plt.show()
```



In [50]:

```
1 #creating stacked column
2 plt.hist([df[df.Genre=="Action"].Budgetmillion, df[df.Genre=="Comedy"].Budgetmillion])
3
4 plt.show()
```



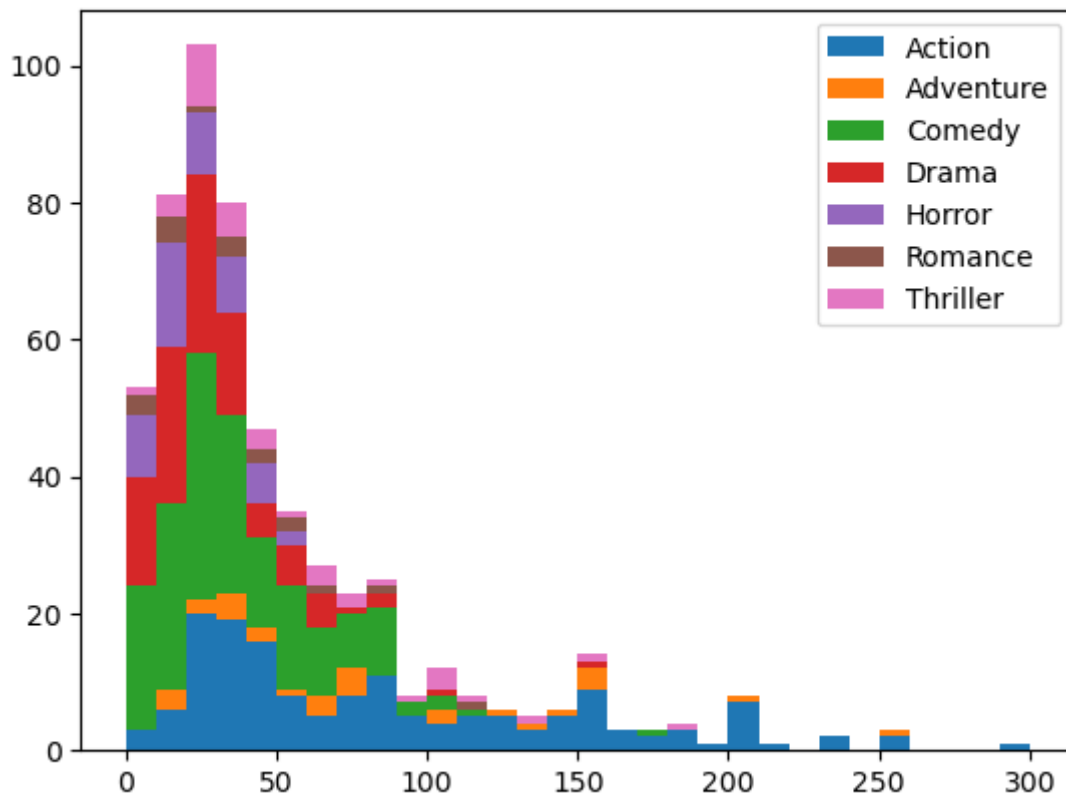
In [51]:

```
1 # so its repetitive task to see all genre distribution in stack
2 #we can automate this task by using loop
3 for gen in df.Genre.cat.categories:
4     print(gen)
```

Action  
Adventure  
Comedy  
Drama  
Horror  
Romance  
Thriller

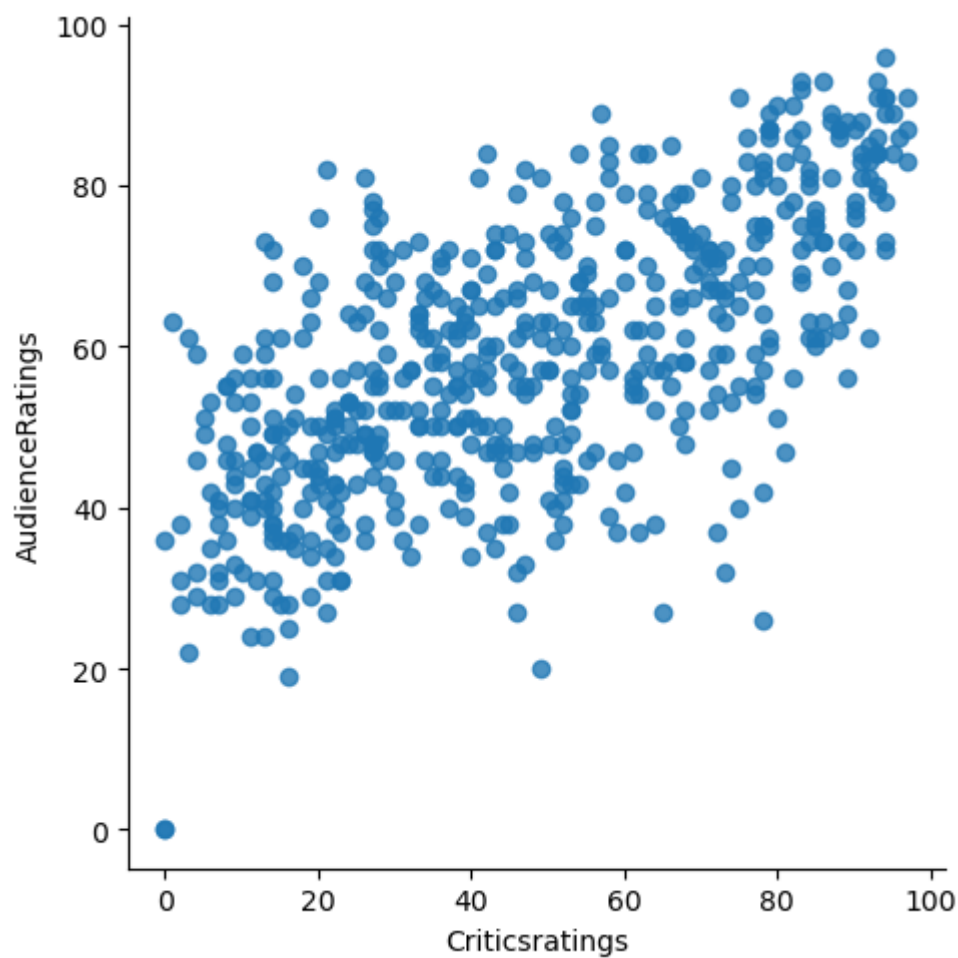
In [65]:

```
1 list1 = list()
2 mylabels=list()
3 for gen in df.Genre.cat.categories:
4     list1.append(df[df.Genre == gen].Budgetmillion)
5     mylabels.append(gen)
6 plt.hist(list1,bins=30,stacked=True,rwidth=1,label=mylabels)
7 plt.legend()
8 plt.show()
```



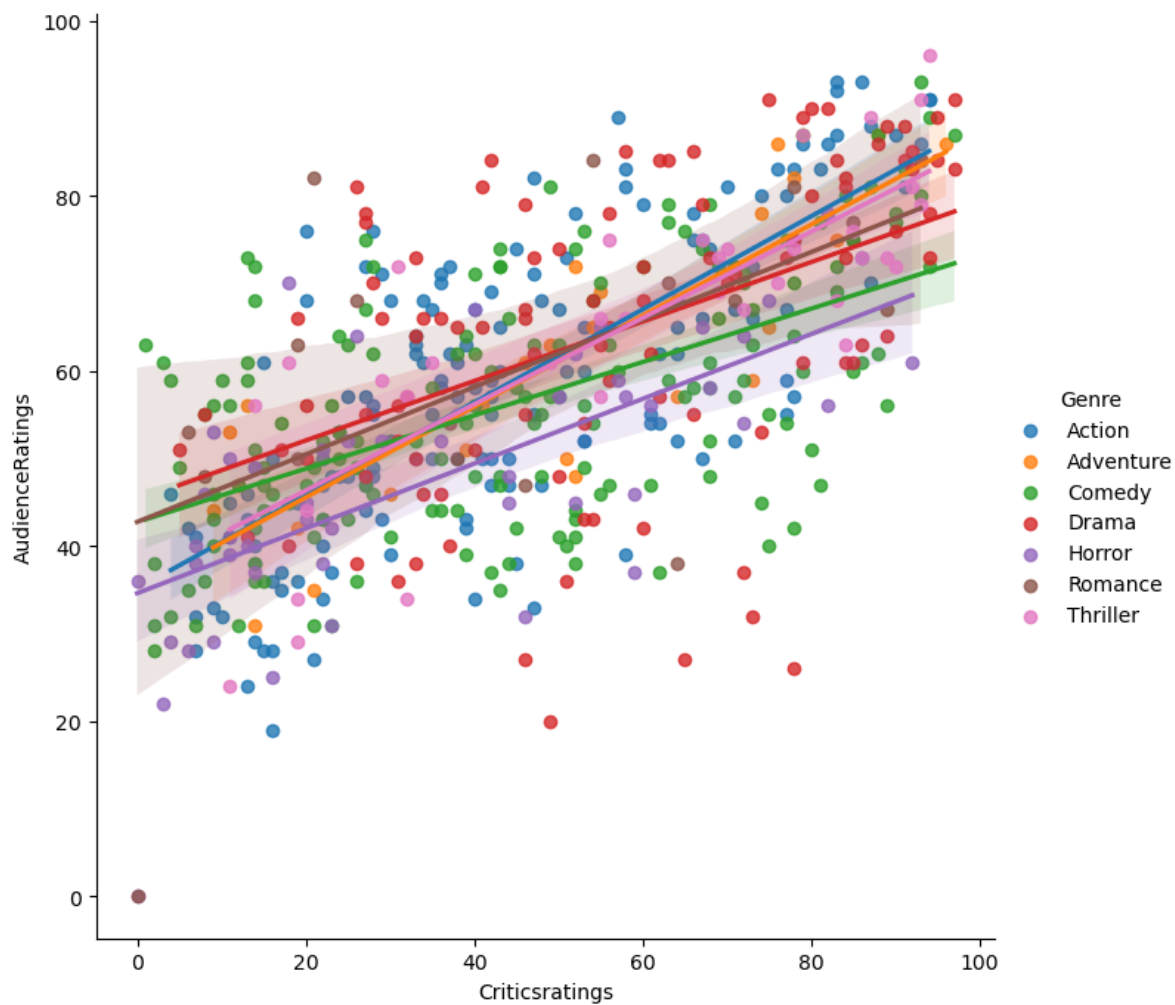
In [66]:

```
1 #KDE plot, code like histplot  
2 vis1=sns.lmplot(data=df, x="Criticsratings", y="AudienceRatings", fit_reg=False)
```



In [68]:

```
1 vis1=sns.lmplot(data=df, x="Criticsratings",y="AudienceRatings",fit_reg=True,hue
```

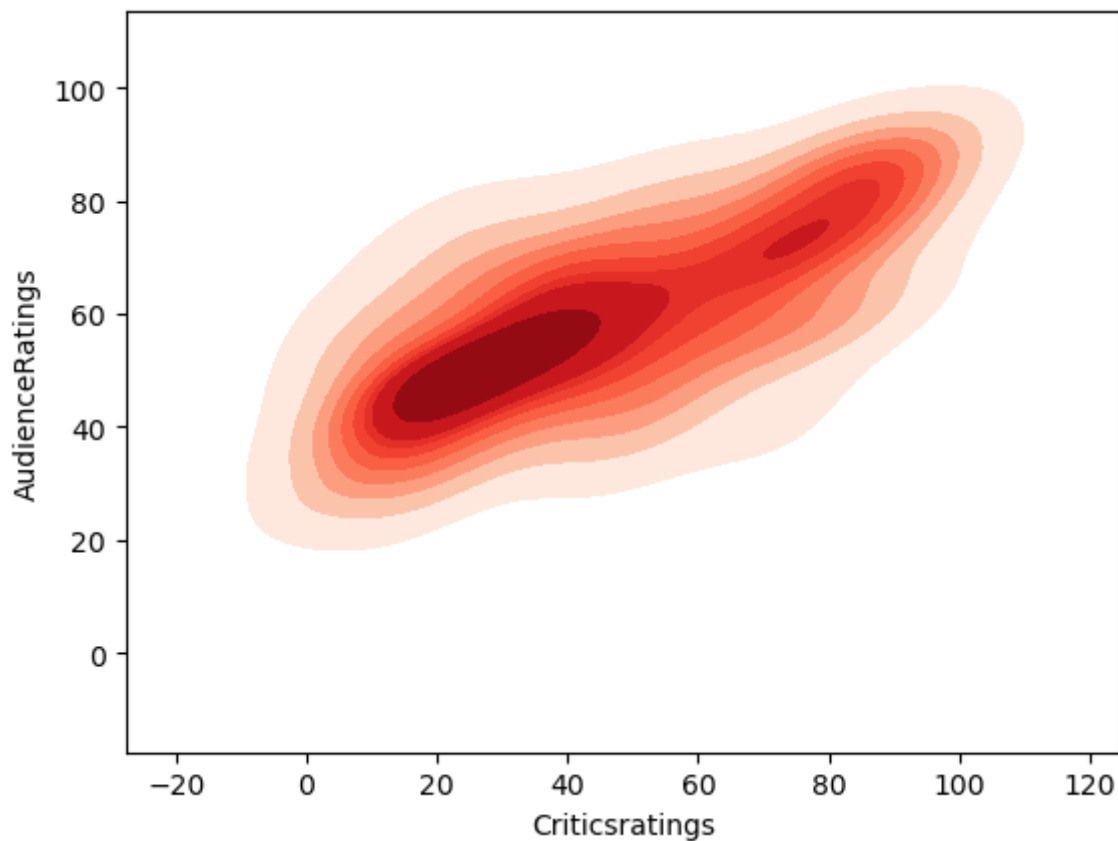


In [76]:

```
1 #KDE
2 import seaborn as sns
3
4 vis2 = sns.kdeplot(data=df, x="Criticsratings", y="AudienceRatings", shade=True,
5 vis2
```

Out[76]:

<Axes: xlabel='Criticsratings', ylabel='AudienceRatings'>

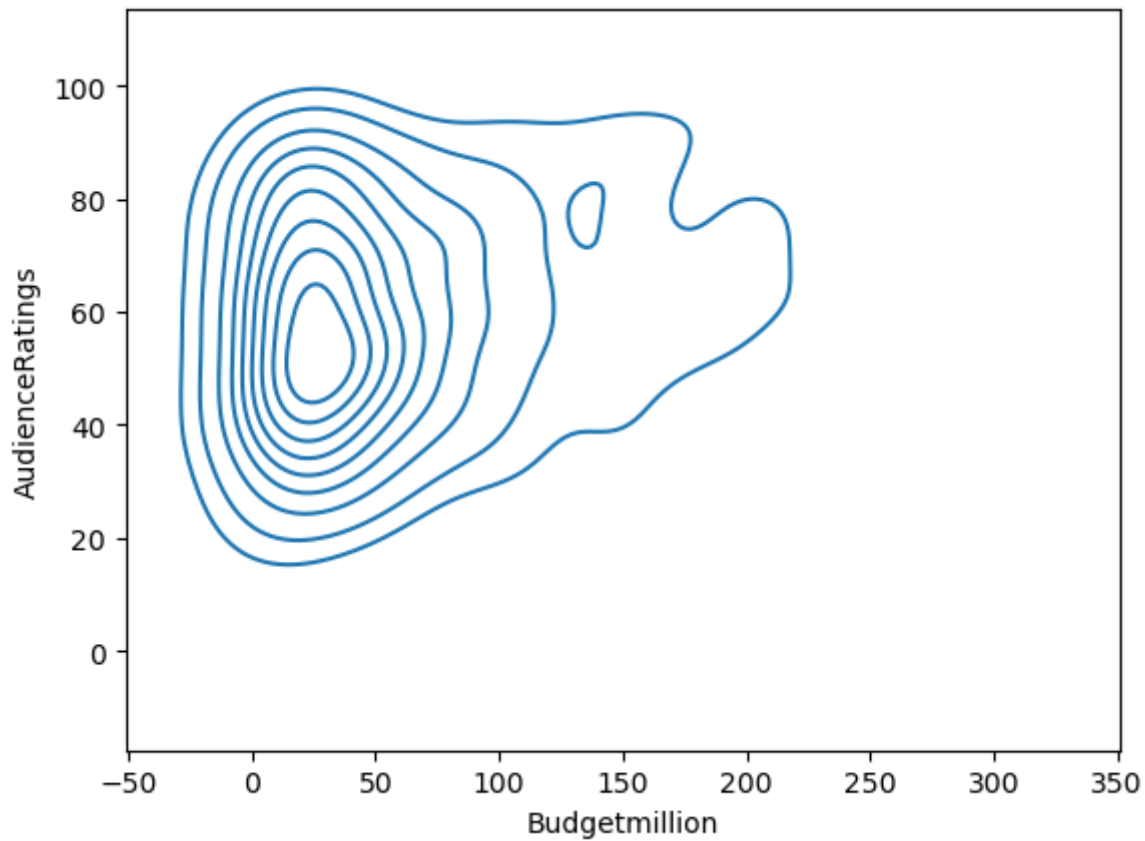


insight:its show where datapoint density is distributed



In [78]:

```
1 #Subplot in python
2 vis3 = sns.kdeplot(data=df, x="Budgetmillion", y="AudienceRatings")
```

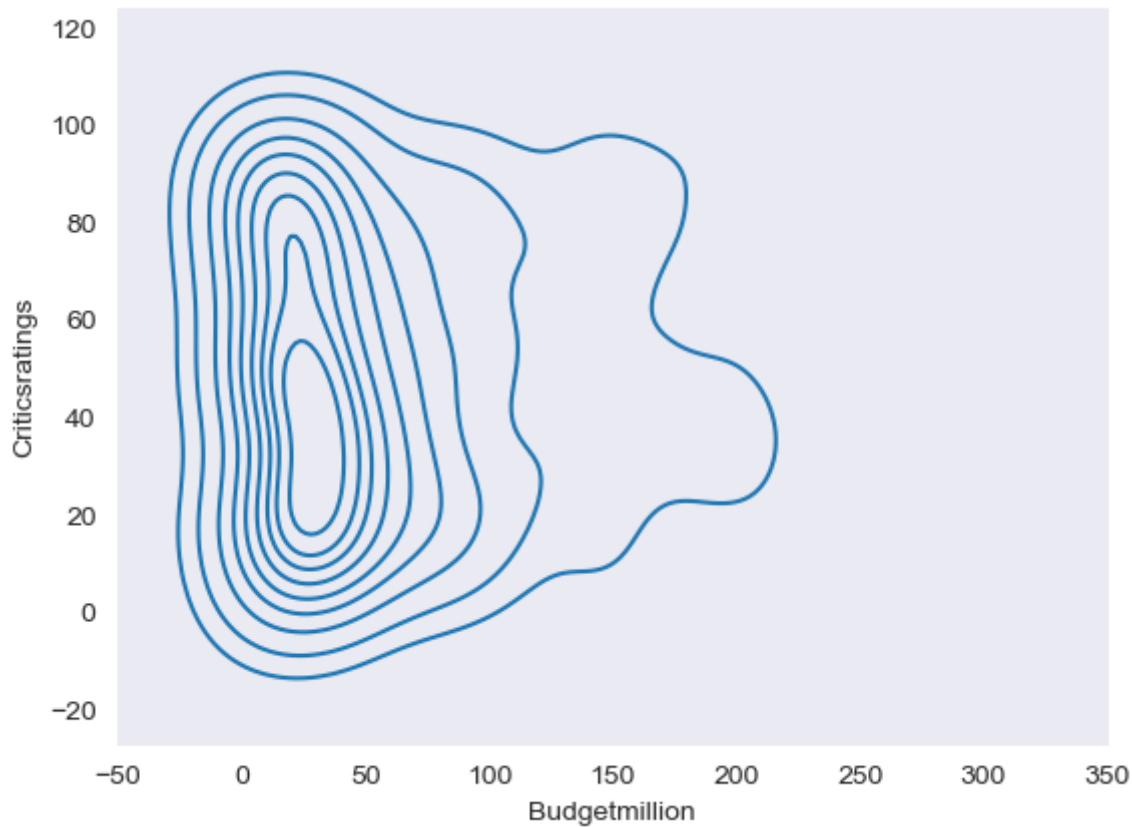


In [ ]:

```
1 insight:how budget affects the audience ratings
```

In [90]:

```
1 sns.set_style("dark")
2 vis4 = sns.kdeplot(data=df, x="Budgetmillion", y="Criticsratings")
```

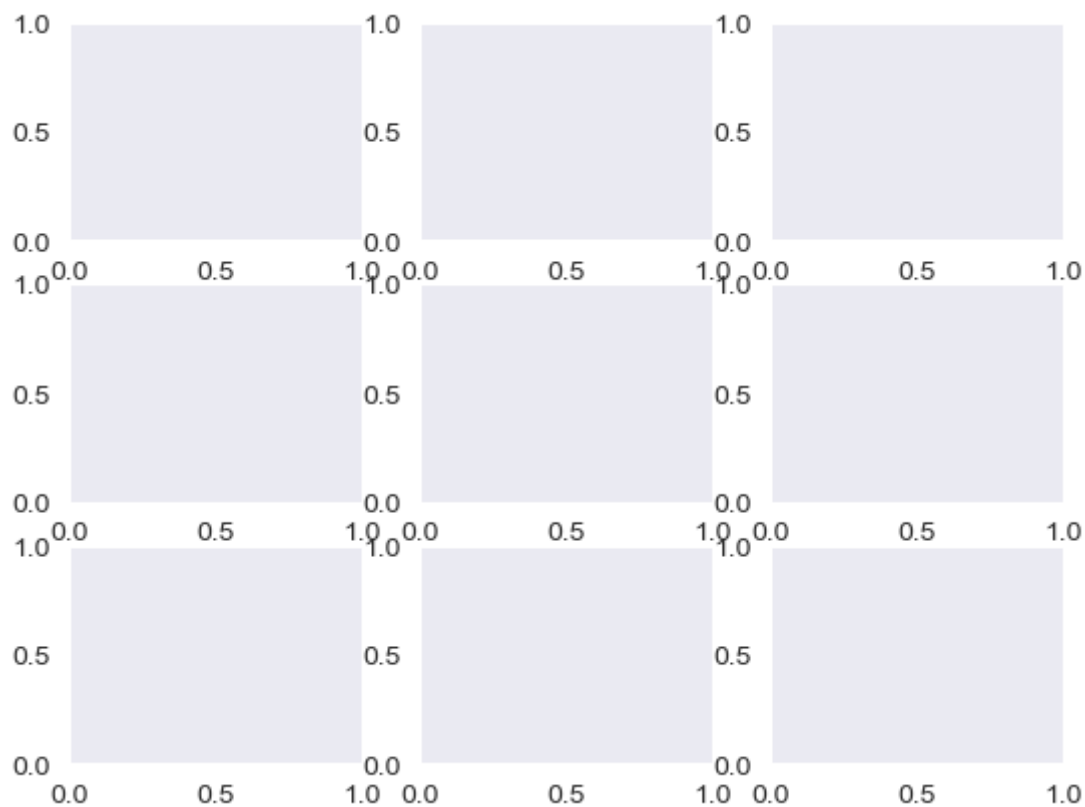


Insight: kind of uniform distribution

## Creating subplot

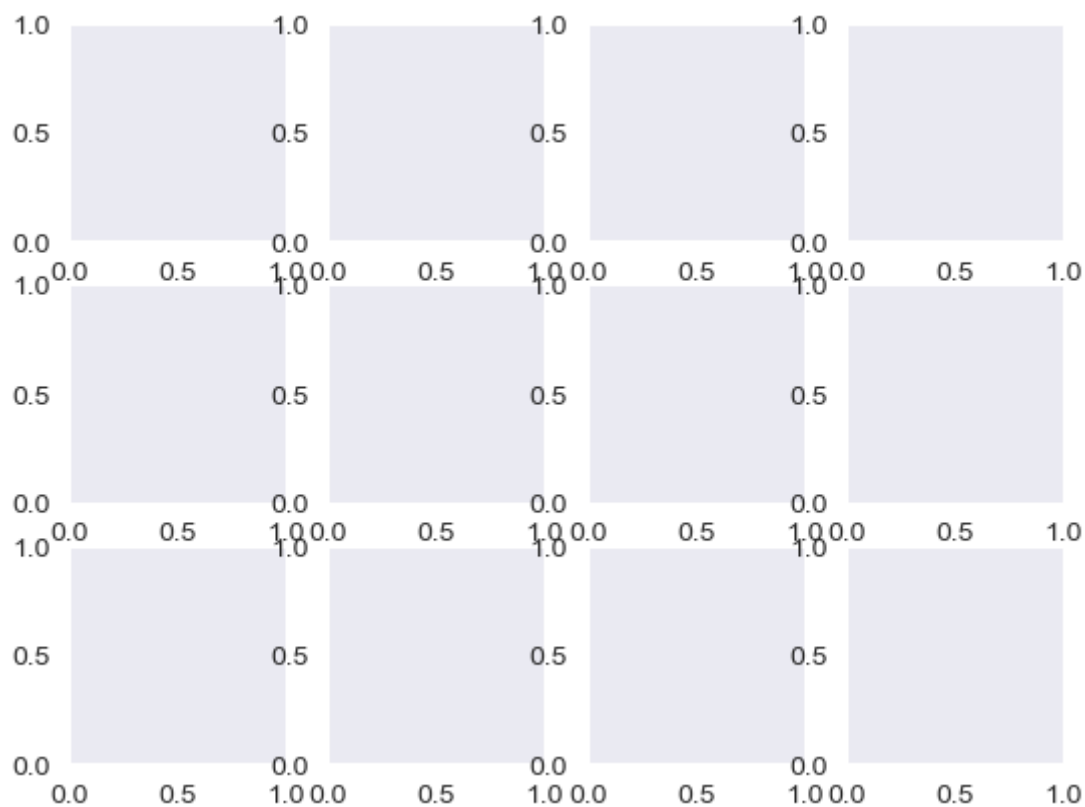
In [85]:

```
1 f,ax=plt.subplots(3,3)
```



In [86]:

```
1 f,ax=plt.subplots(3,4)
```



In [94]:

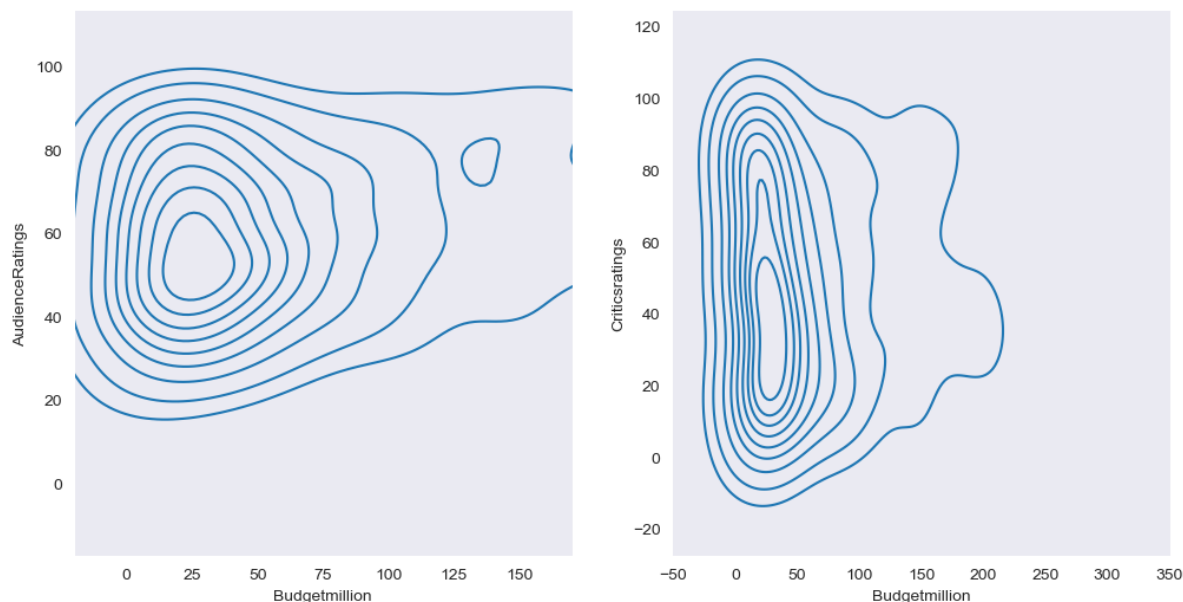
```

1 f,axes=plt.subplots(1,2,figsize=(12,6))
2 vis3 = sns.kdeplot(data=df, x="Budgetmillion", y="AudienceRatings",ax=axes[0])
3 vis4 = sns.kdeplot(data=df, x="Budgetmillion", y="Criticsratings",ax=axes[1])
4 vis3.set(xlim=(-20,170))

```

Out[94]:

[(-20.0, 170.0)]



In [161]:

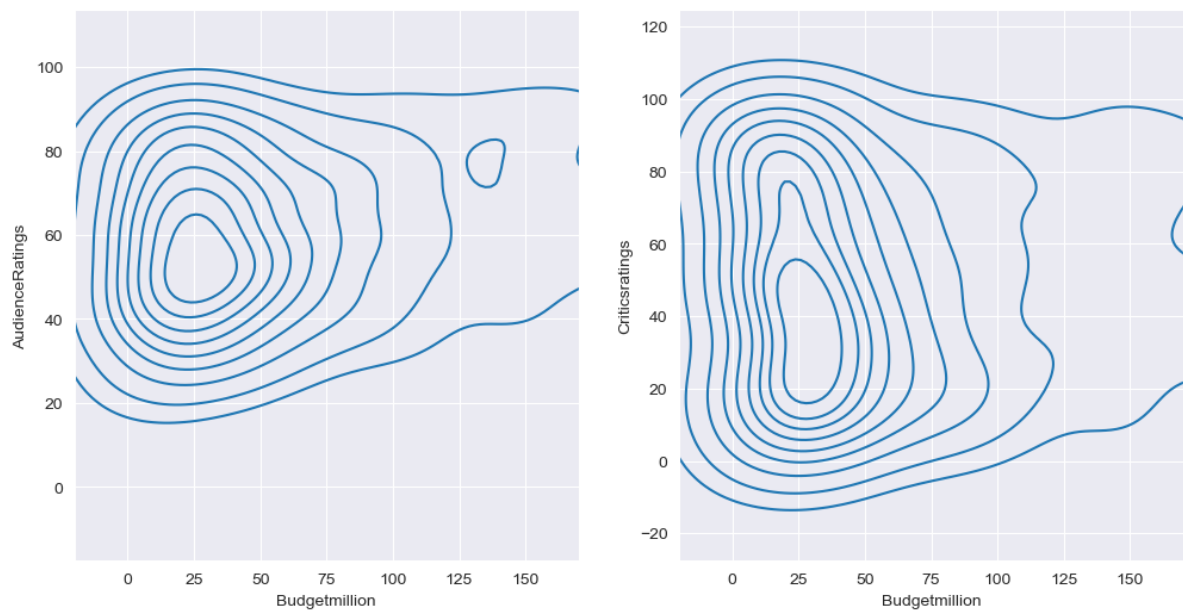
```

1 f,axes=plt.subplots(1,2,figsize=(12,6),sharex=True)
2 vis3 = sns.kdeplot(data=df, x="Budgetmillion", y="AudienceRatings",ax=axes[0])
3 vis4 = sns.kdeplot(data=df, x="Budgetmillion", y="Criticsratings",ax=axes[1])
4 vis3.set(xlim=(-20,170))

```

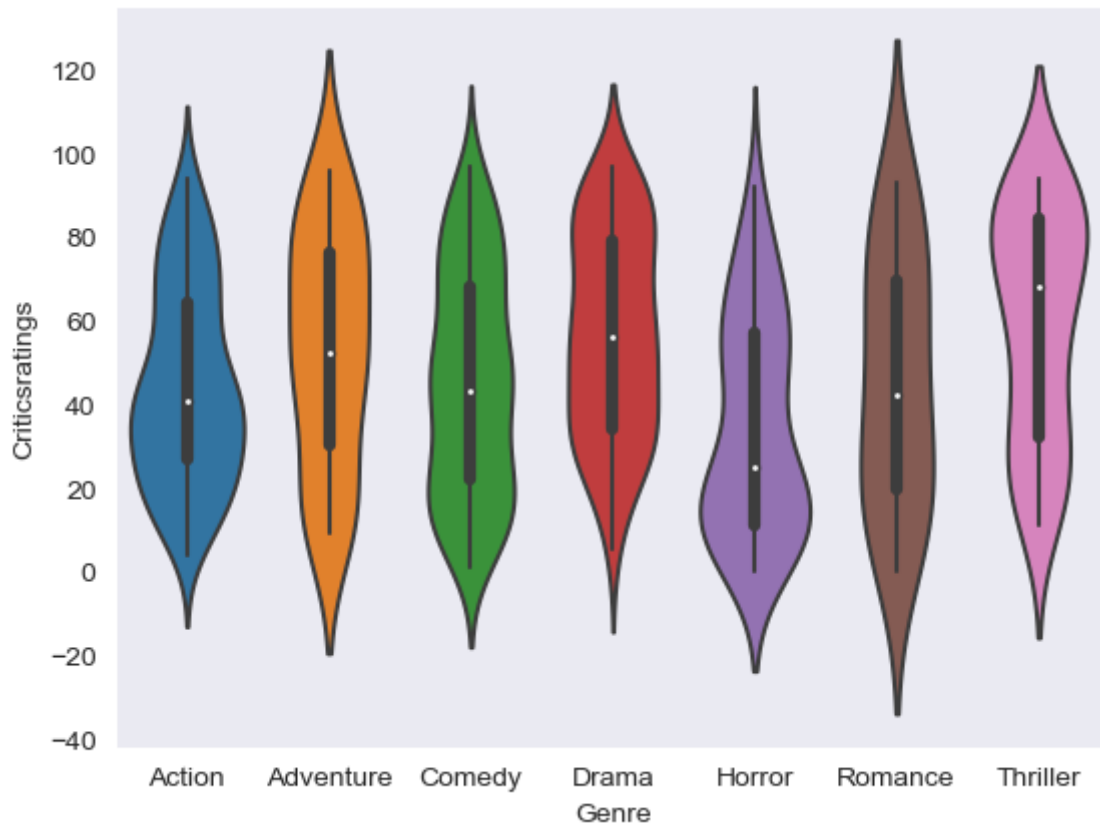
Out[161]:

[(-20.0, 170.0)]



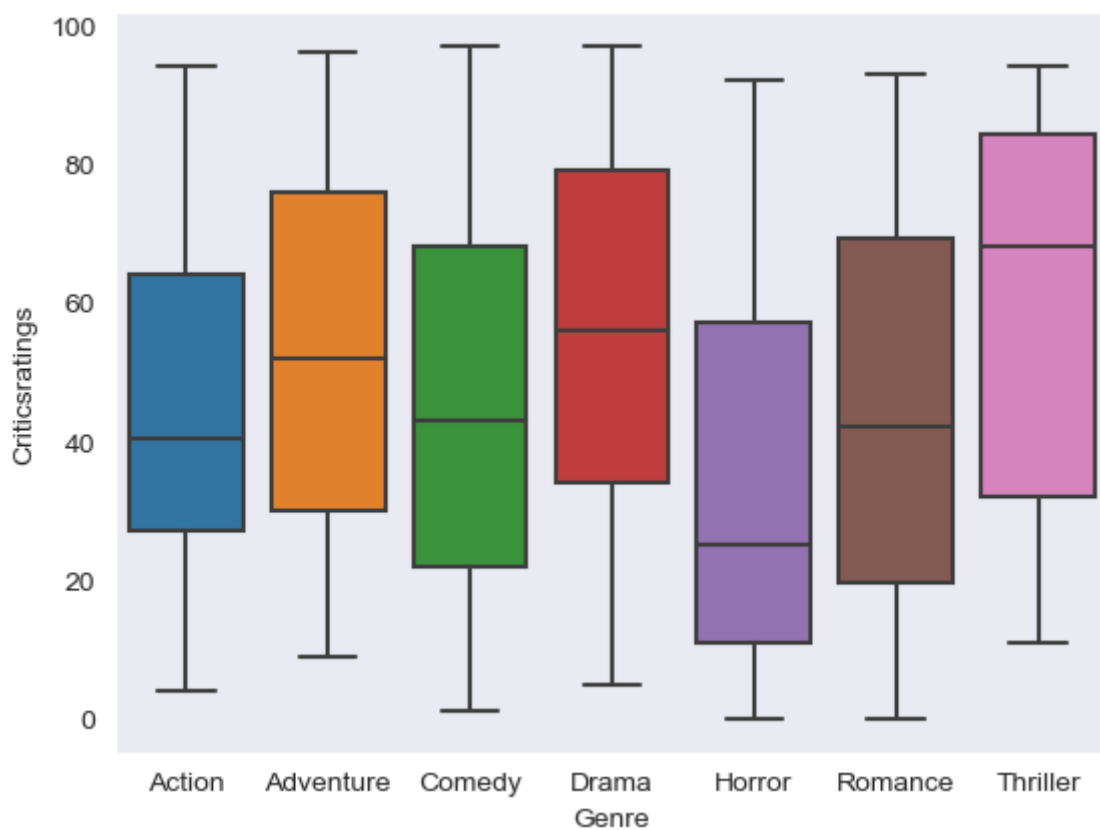
In [96]:

```
1 # violin and box plot
2 v=sns.violinplot(data=df,x="Genre",y="Criticsratings")
```



In [97]:

```
1 b=sns.boxplot(data=df,x="Genre",y="Criticsratings")
```



Insight :its tells hows critic rating distributed across genre, first quartile 2nd and 3rd quartile, , median(chck the highest or lowest median among them) max min and outlier inthis case thriller have high median for critic sratings, so we will build thriller

In [99]:

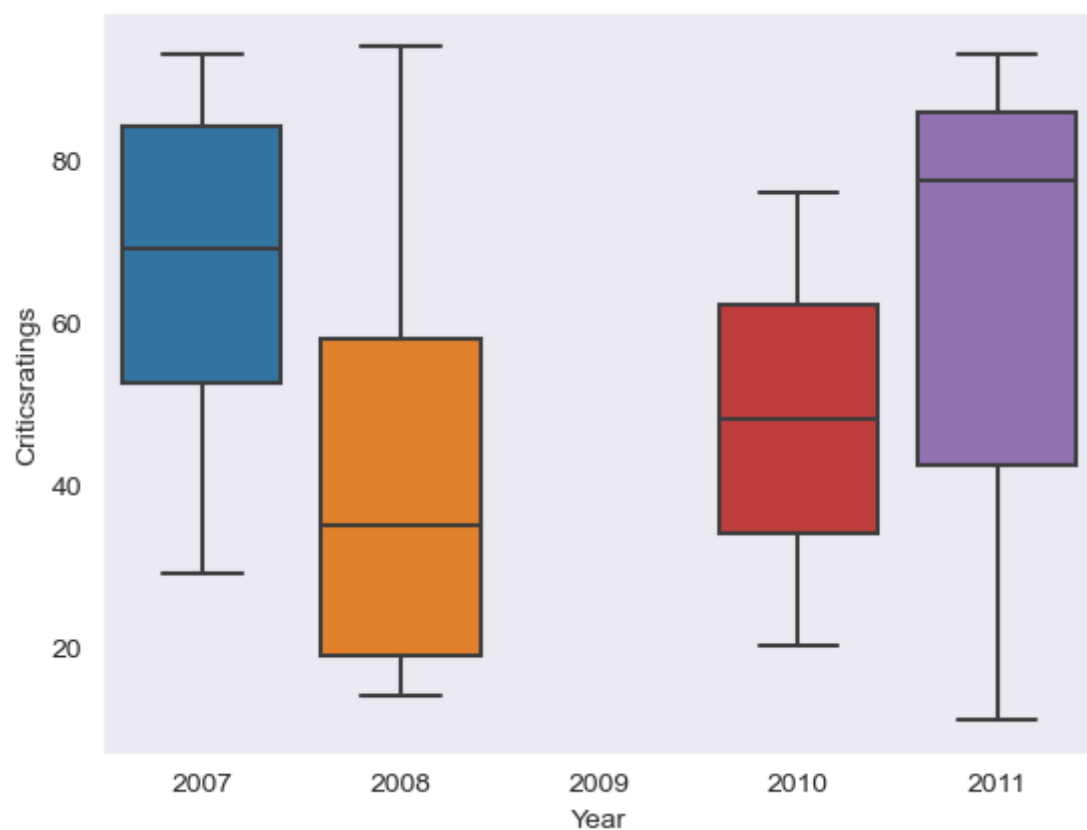
```
1 #if you want to see only one genre, critics ratings abd with different year
2 # first crate a new dta frame for specific genre
3 Thriller=df[df.Genre=="Thriller"]
```

In [100]:

```
1 b=sns.boxplot(data=Thriller,x="Year",y="Criticsratings")
2 b
```

Out[100]:

<Axes: xlabel='Year', ylabel='Criticsratings'>

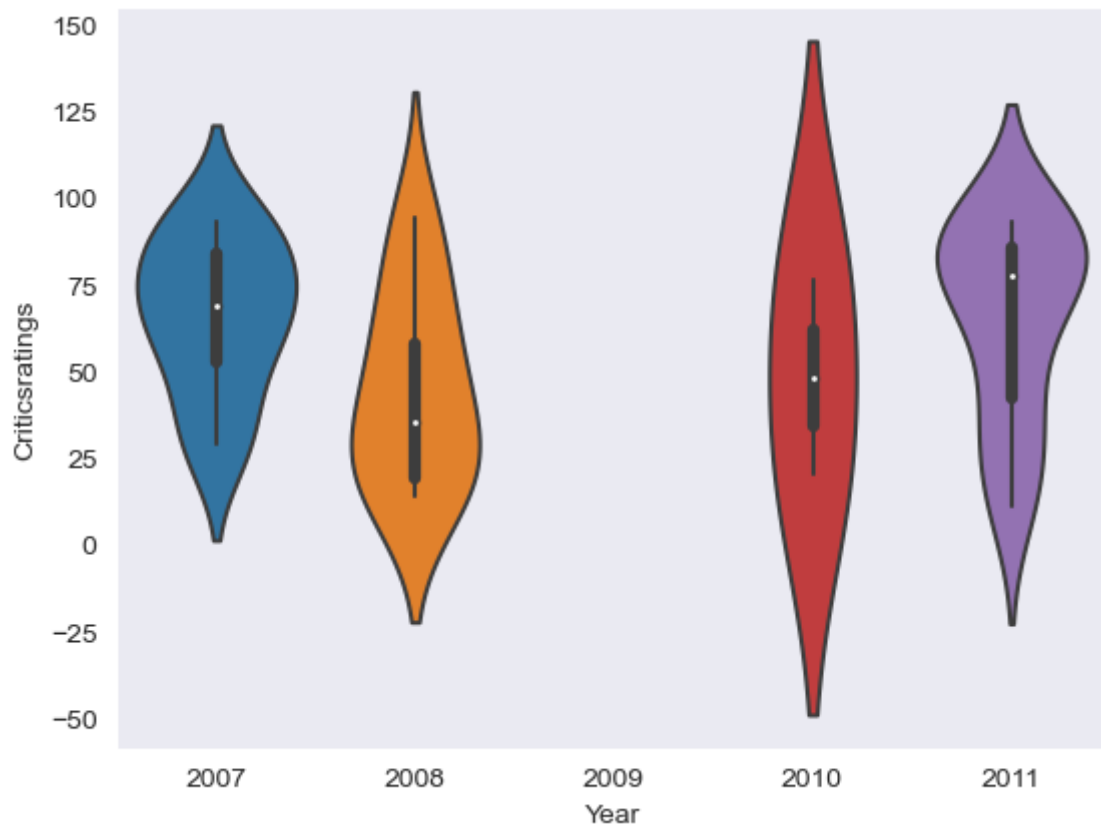


In [101]:

```
1 vl=sns.violinplot(data=Thriller,x="Year",y="Criticsratings")
2 vl
```

Out[101]:

&lt;Axes: xlabel='Year', ylabel='Criticsratings'&gt;



In [ ]:

```
1 Insight: Violinplot can show you something better than boxplot by shape of box
```

#creating facedgrid #how to create multiple chart for separate genre, why ?if you see lmplo, we can identify separately datapoint for individual genre

In [110]:

```
1 g=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
```



in here, you basically controles the visualizaion using some sort rules

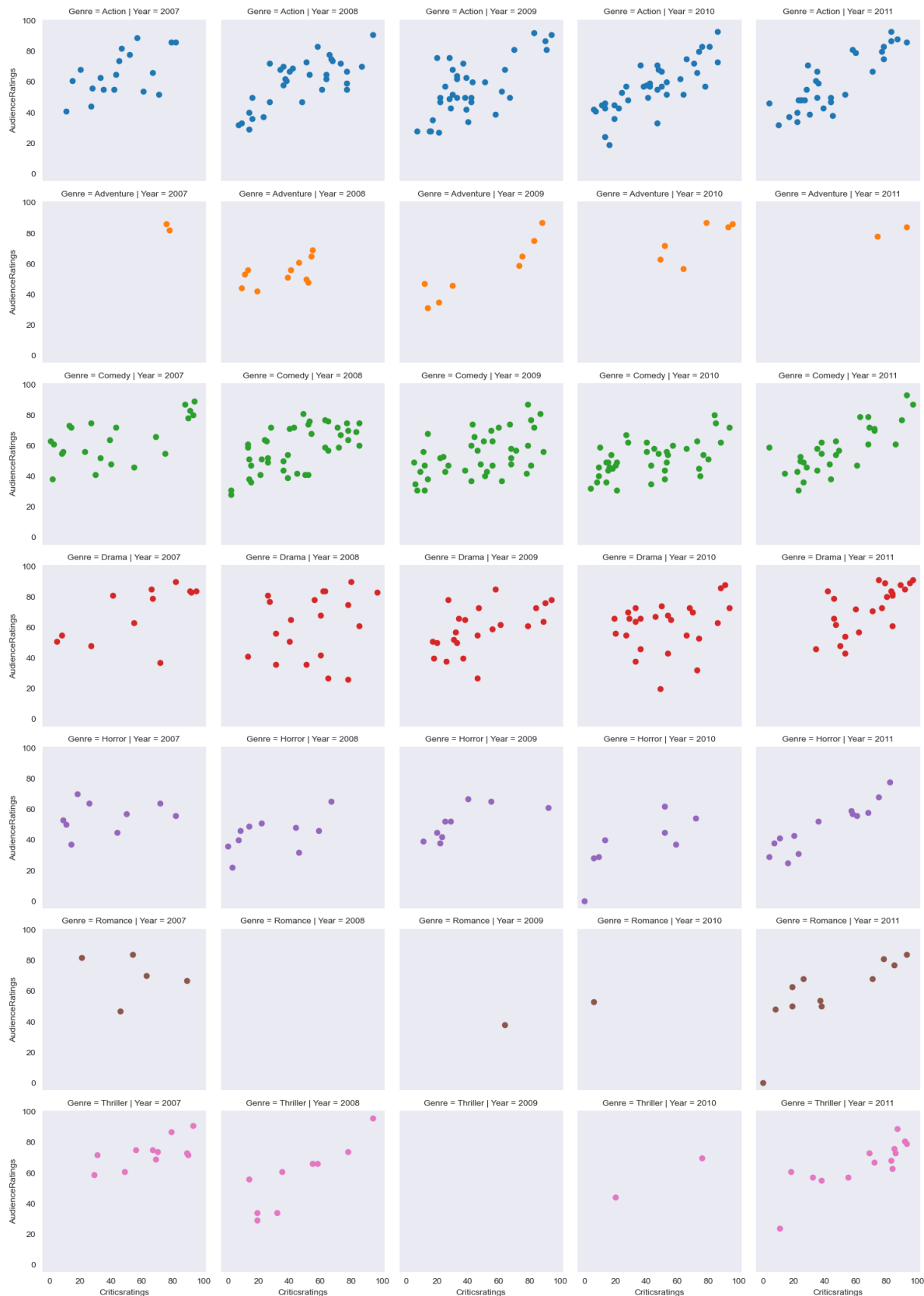


In [115]:

```

1 #now we are going to insert plot into facedgrid by map function
2 g=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
3 g=g.map(plt.scatter,"Criticsratings","AudienceRatings")
4

```



In [118]:

```

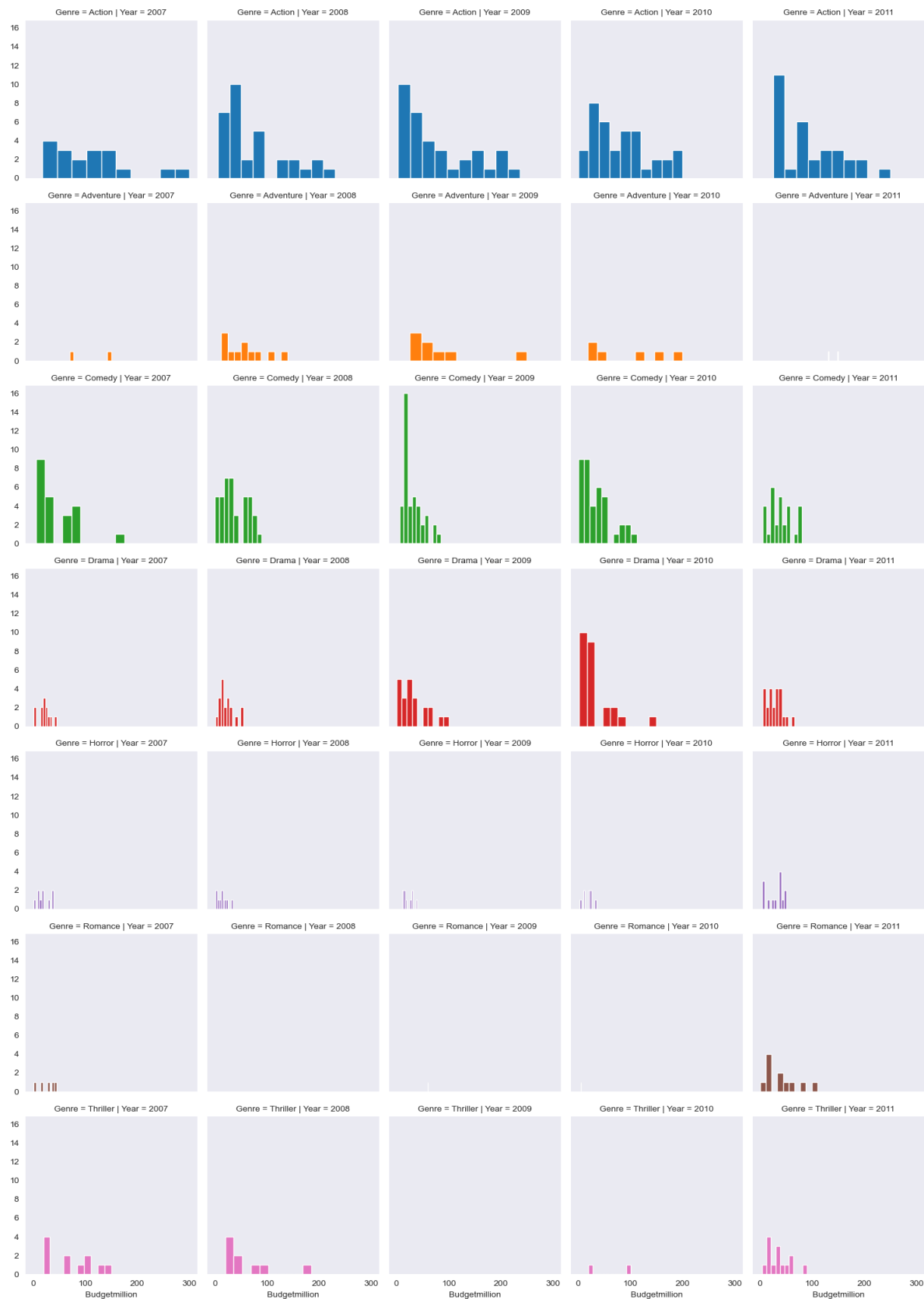
1 # you can do this for other plot also
2 #do for histplot
3 h=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
4 h= h.map(plt.hist,"Criticsratings")
5

```



In [120]:

```
1 h=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
2 h=h.map(plt.hist,"Budgetmillion")
```

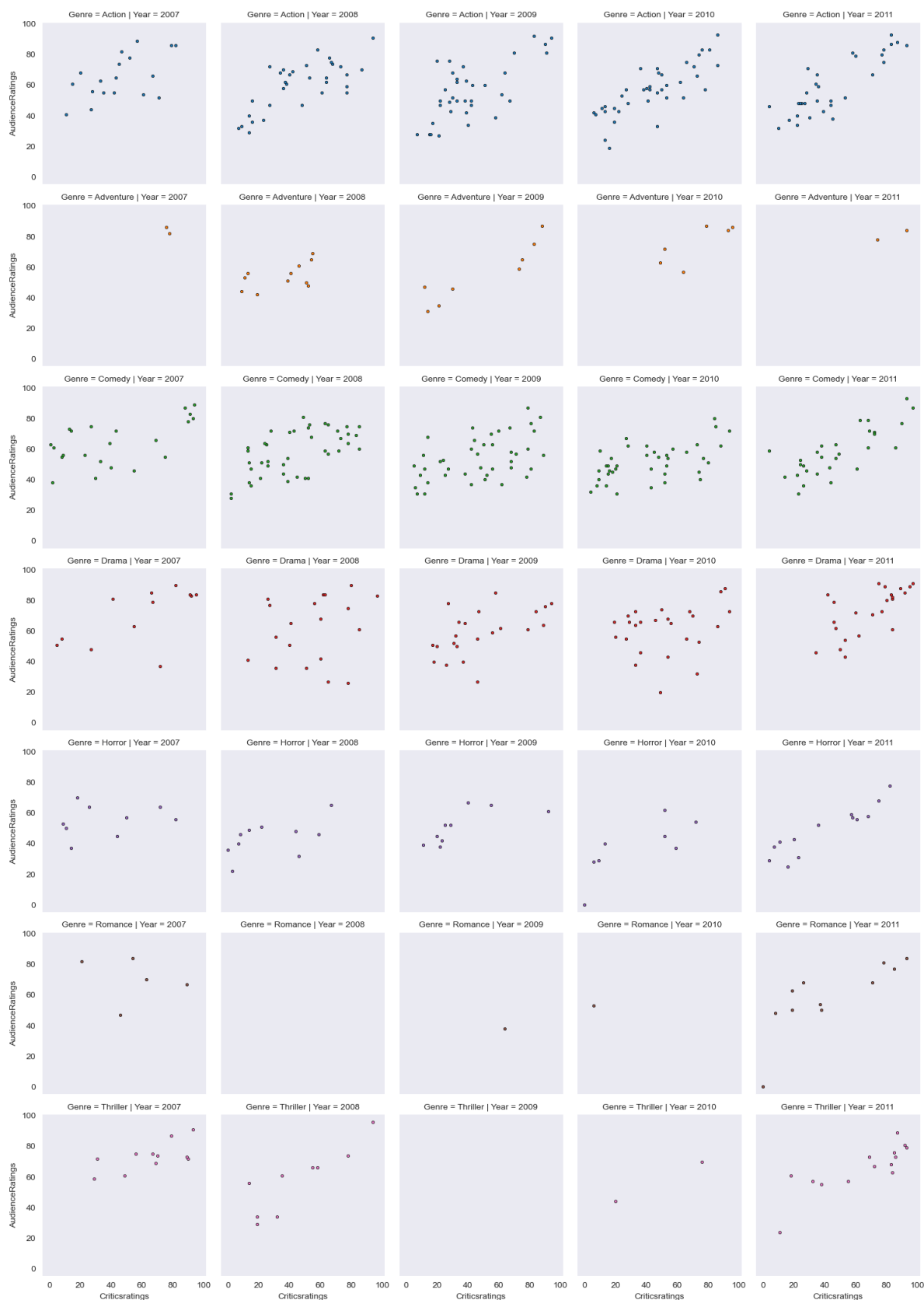


In [123]:

```

1 #inserting third argument , like contoure of data point size
2
3 g=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
4 kws=dict(s=10,linewidth=.5,edgecolor="black")
5 g=g.map(plt.scatter,"Criticsratings","AudienceRatings",**kws)
6

```



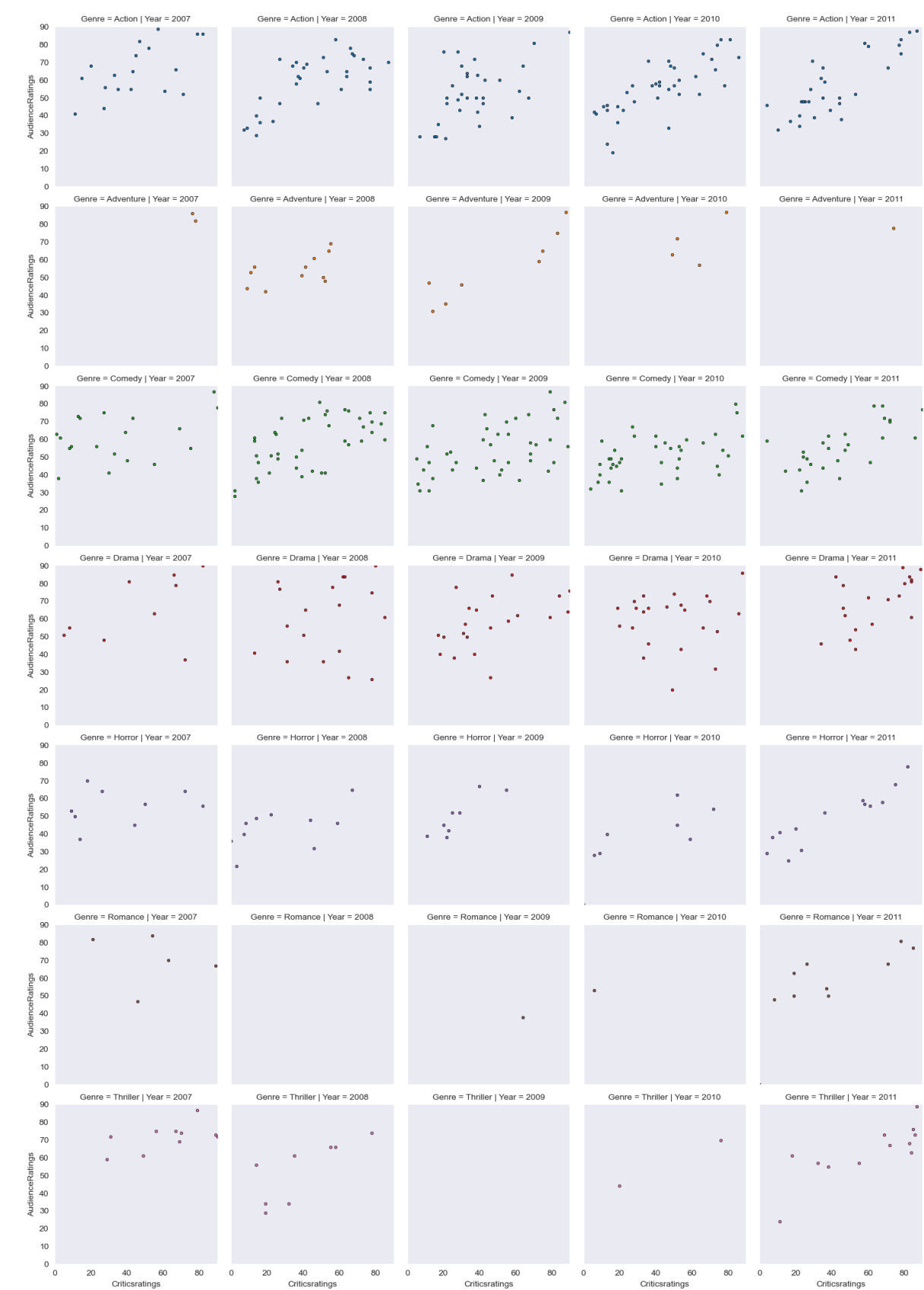
insight: tells how audience and critics rating for each genre throught the year

In [127]:

```
1 #optimizing cordinal and diagonal to bettter understans, if you see cardinal (x
2 #and by diagonal you can distinguished well critics ratings and audience ratings
3 # for this cardinal, set the limit for x and y axes
4 g=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
5 kws=dict(s=10,linewidth=.5,edgecolor="black")
6 g=g.map(plt.scatter,"Criticsratings","AudienceRatings",**kws)
7 g.set(xlim=(0,90),ylim=(0,90))
```

Out[127]:

<seaborn.axisgrid.FacetGrid at 0x7faa3a6ef220>

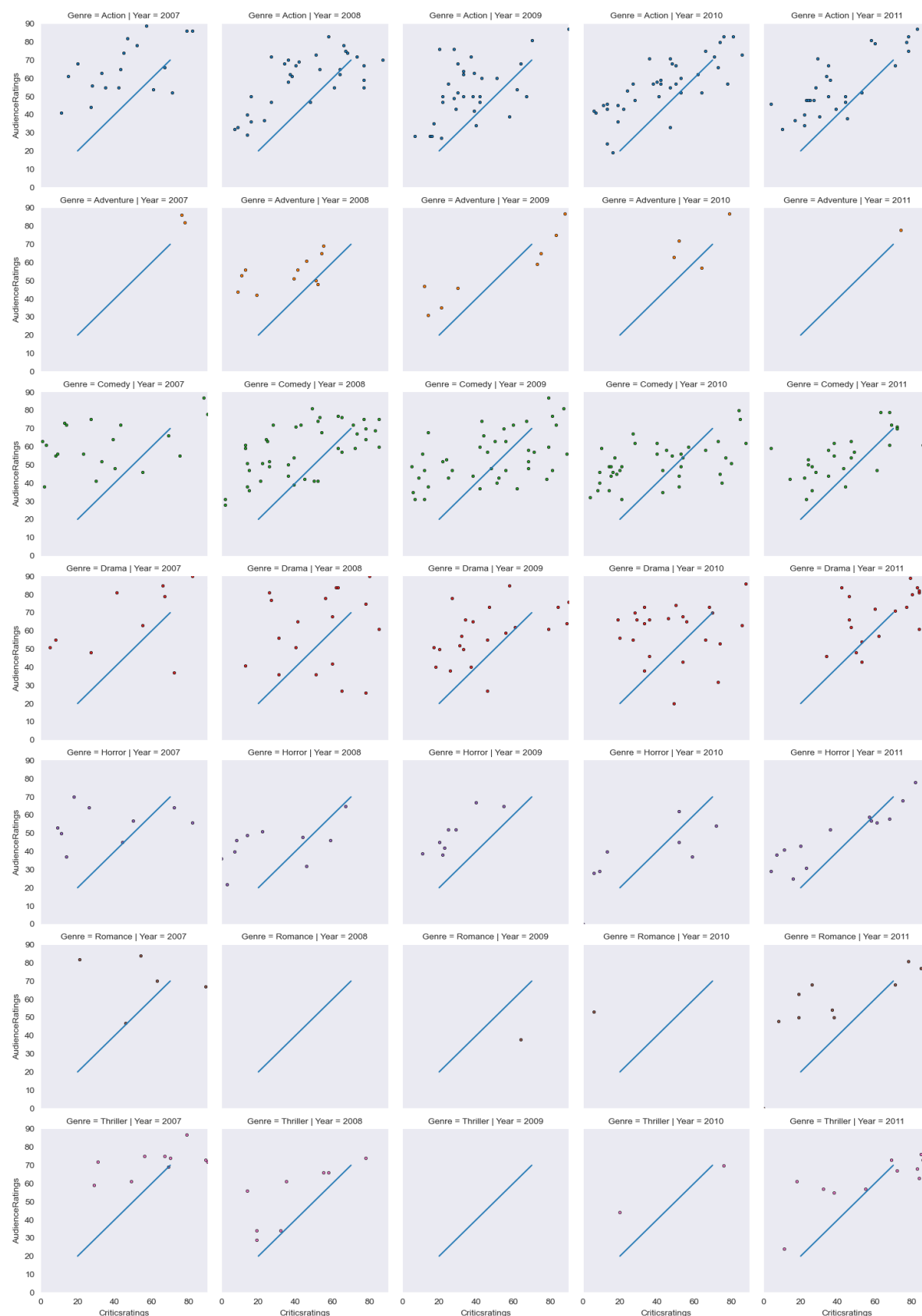


In [129]:

```

1 #for insert diagonal , you have to loop in ax andc axes as array
2 # thhen define the range from what poits draw the line
3 g=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
4 kws=dict(s=10,linewidth=.5,edgecolor="black")
5 g=g.map(plt.scatter,"Criticsratings","AudienceRatings",**kws)
6 g.set(xlim=(0,90),ylim=(0,90))
7 for ax in g.axes.flat:
8     ax.plot((20,70),(20,70))

```



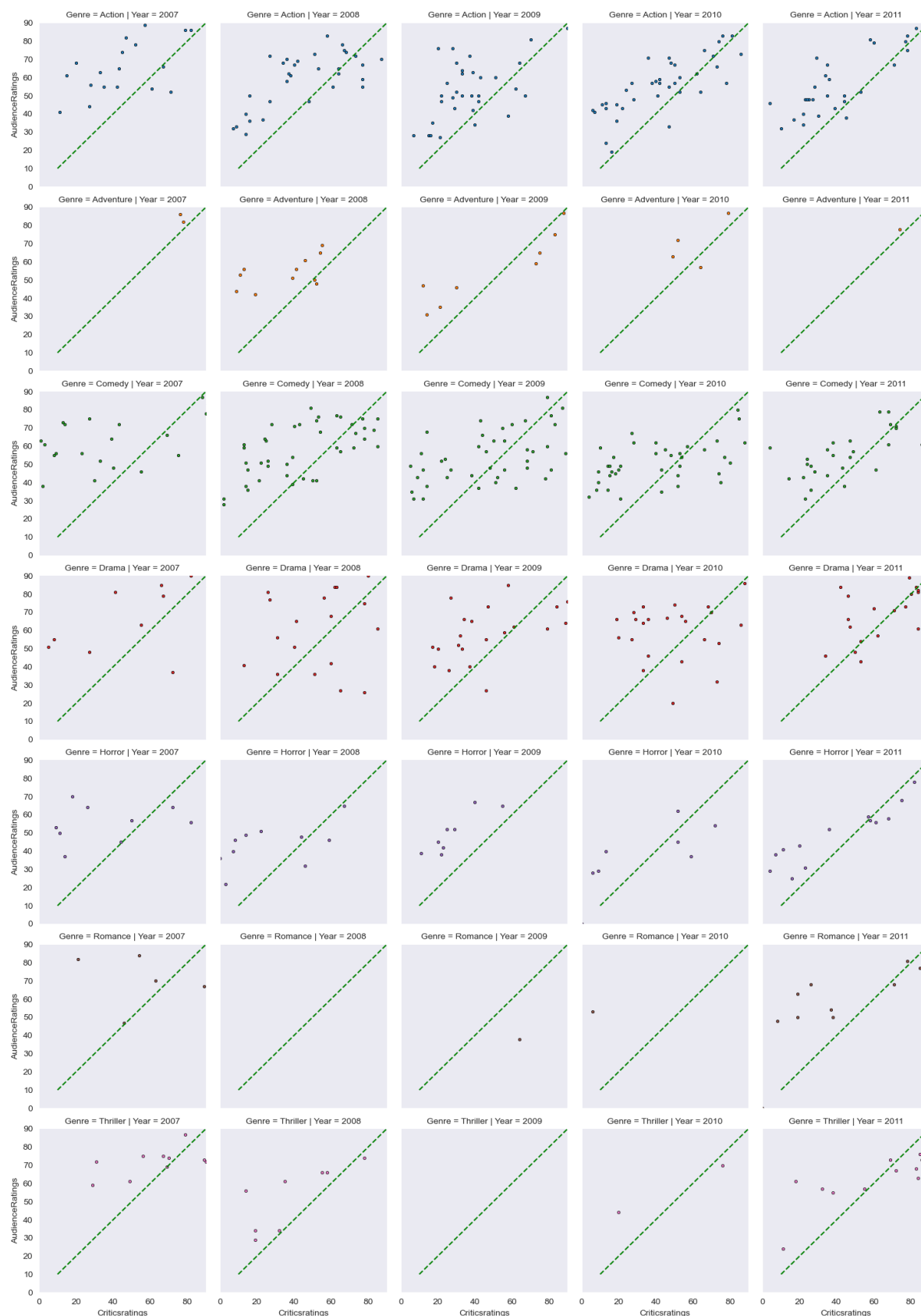


In [131]:

```

1 g=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
2 kws=dict(s=10,linewidth=.5,edgecolor="black")
3 g=g.map(plt.scatter,"Criticsratings","AudienceRatings",**kws)
4 g.set(xlim=(0,90),ylim=(0,90))
5 for ax in g.axes.flat:
6     ax.plot((10,100),(10,100),c="green",ls="--")

```



Insight:so this diagonal helps to detect where audince rating is more (if upper the diagonal more point) or vice versa

In [133]:

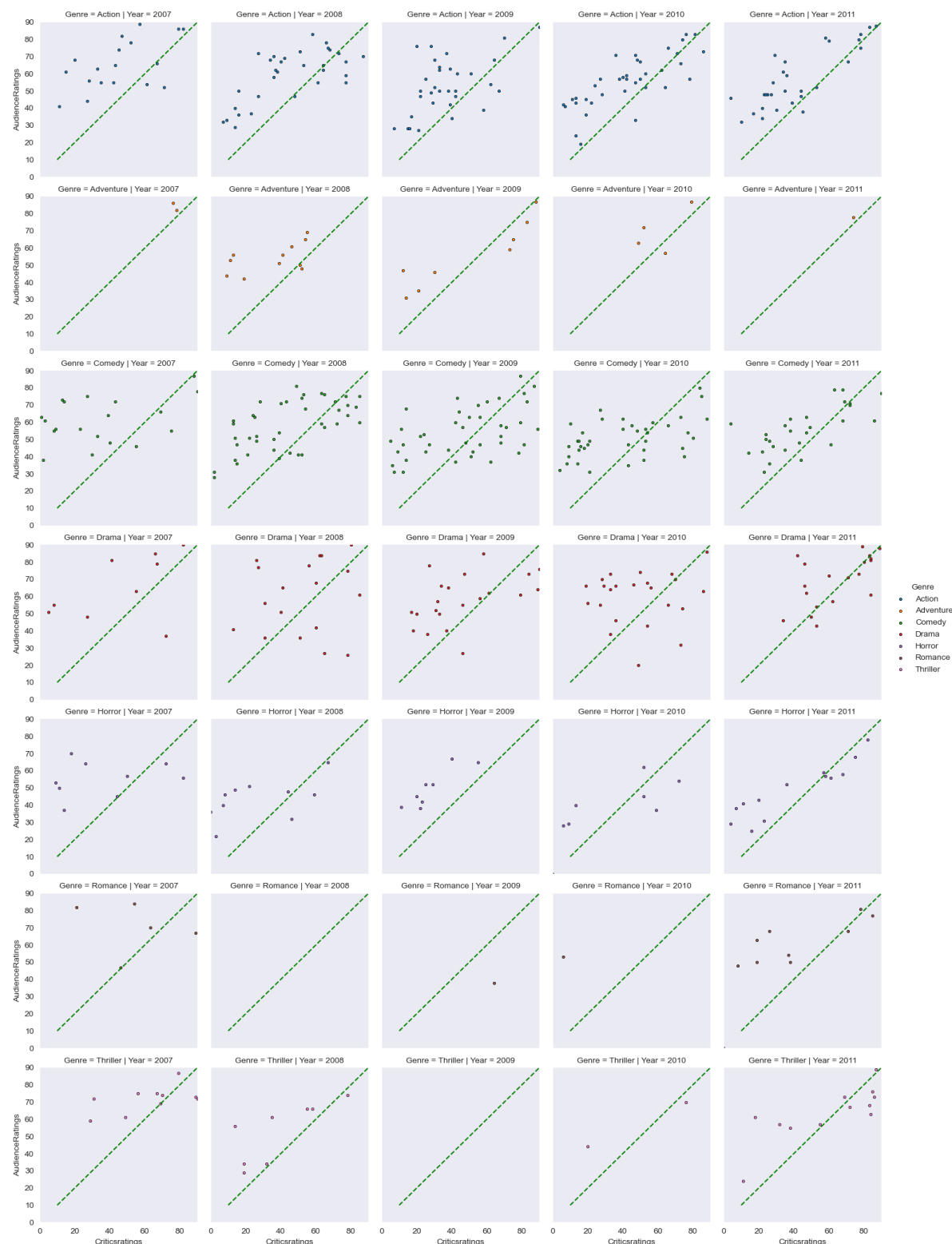
```

1 #adding legend
2 g=sns.FacetGrid(df,row="Genre",col="Year",hue="Genre")
3 kws=dict(s=10,linewidth=.5,edgecolor="black")
4 g=g.map(plt.scatter,"Criticsratings","AudienceRatings",**kws)
5 g.set(xlim=(0,90),ylim=(0,90))
6 for ax in g.axes.flat:
7     ax.plot((10,100),(10,100),c="green",ls="--")
8 g.add_legend()

```

Out[133]:

&lt;seaborn.axisgrid.FacetGrid at 0x7faa3d1818a0&gt;

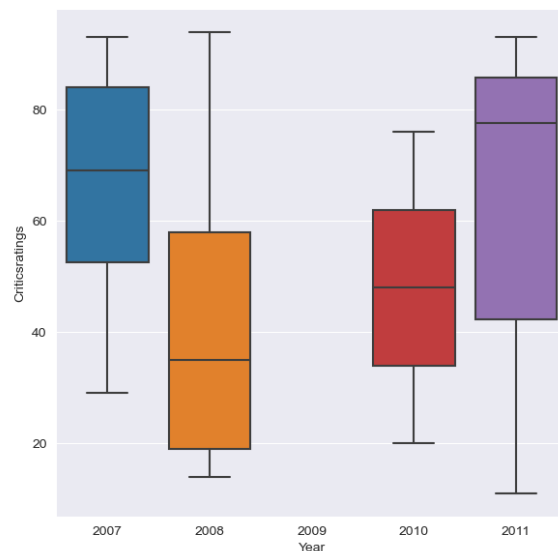
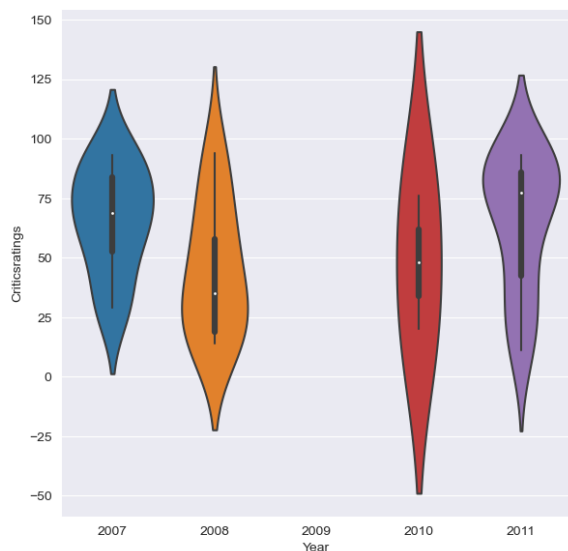
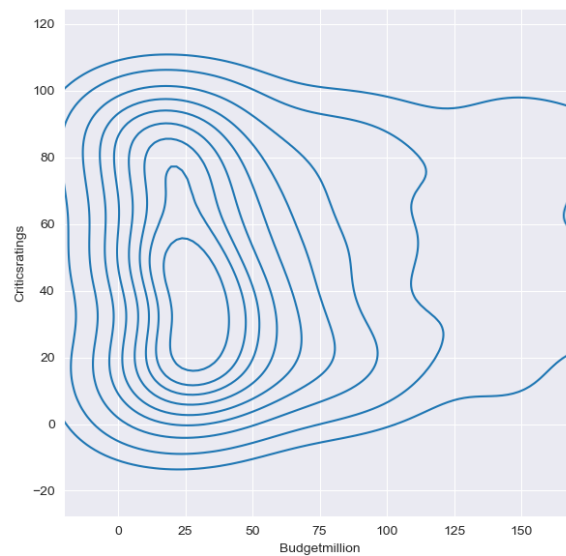
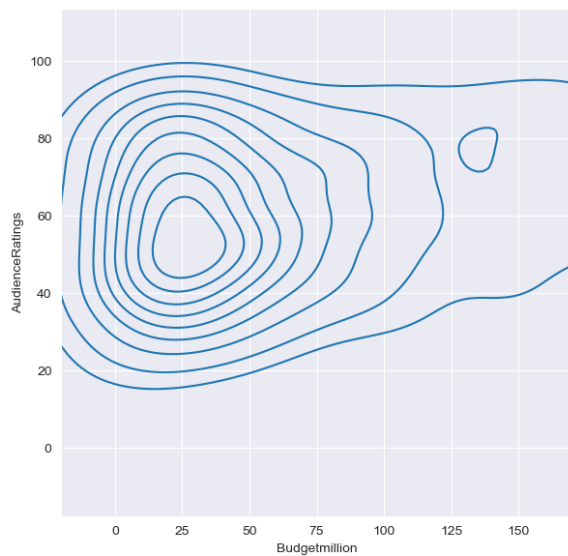


# Dashboard

Dashboard is combination of chart, a figure with multiple chart, we will use subplot function to create dashboard

In [166]:

```
1 sns.set_style("darkgrid")
2 f, axes = plt.subplots(2, 2, figsize=(15, 15))
3 vis3 = sns.kdeplot(data=df, x="Budgetmillion", y="AudienceRatings", ax=axes[0, 0])
4 vis4 = sns.kdeplot(data=df, x="Budgetmillion", y="Criticsratings", ax=axes[0, 1])
5 vis3.set(xlim=(-20, 170))
6 vis4.set(xlim=(-20, 170))
7 vl=sns.violinplot(data=Thriller, x="Year", y="Criticsratings", ax=axes[1, 0])
8 b=sns.boxplot(data=Thriller, x="Year", y="Criticsratings", ax=axes[1, 1])
```

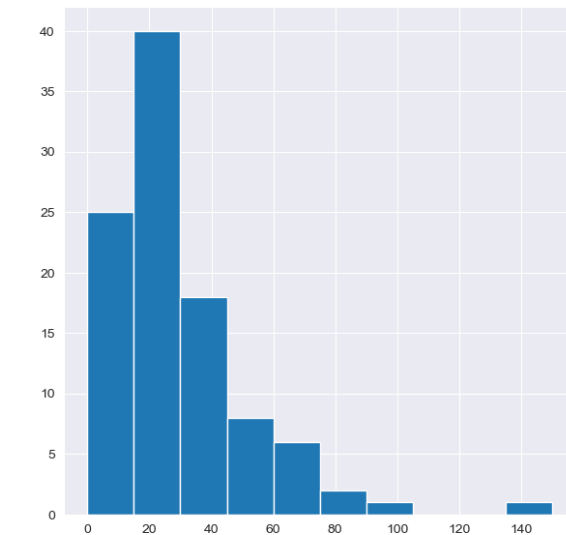
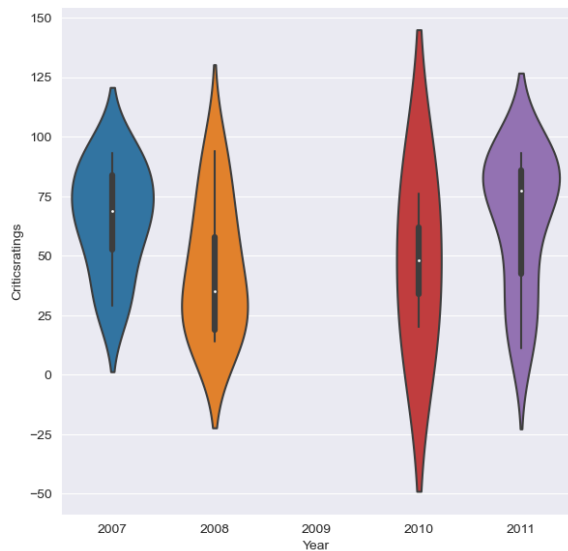
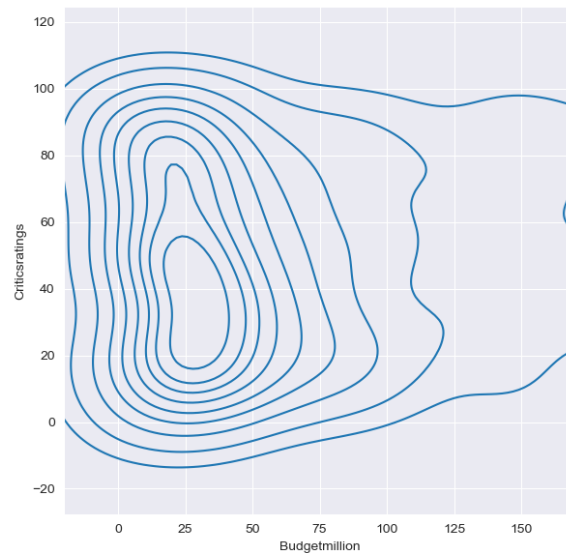
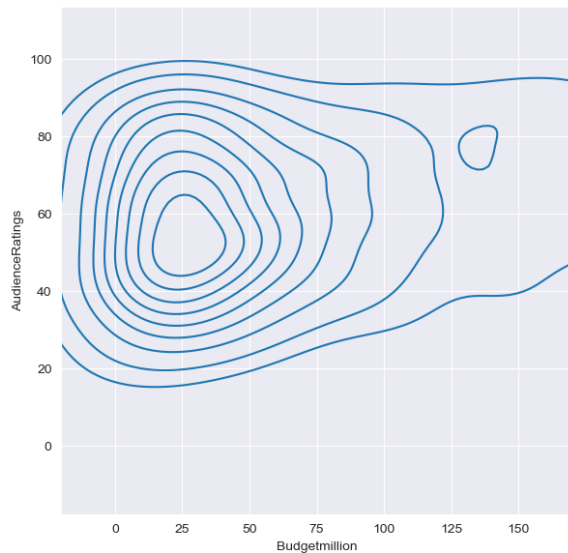


In [169]:

```

1 #if you want to add std.plot to the seaborn plot , you have to apply a different
2 sns.set_style("darkgrid")
3 f,axes=plt.subplots(2,2,figsize=(15,15))
4 vis3 = sns.kdeplot(data=df, x="Budgetmillion", y="AudienceRatings",ax=axes[0,0])
5 vis4 = sns.kdeplot(data=df, x="Budgetmillion", y="Criticsratings",ax=axes[0,1])
6 vis3.set(xlim=(-20,170))
7 vis4.set(xlim=(-20,170))
8 vl=sns.violinplot(data=Thriller,x="Year",y="Criticsratings",ax=axes[1,0])
9 axes[1,1].hist(df[df.Genre=="Drama"].Budgetmillion)
10 plt.show()

```

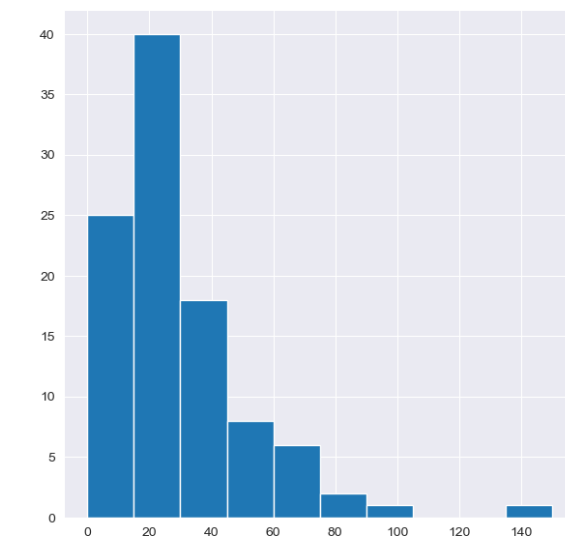
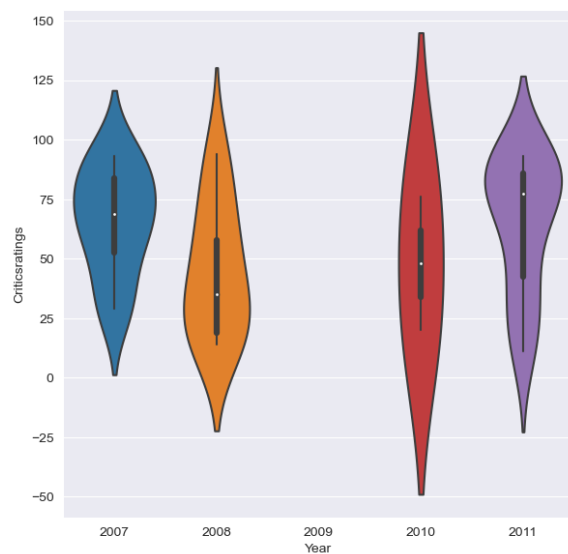
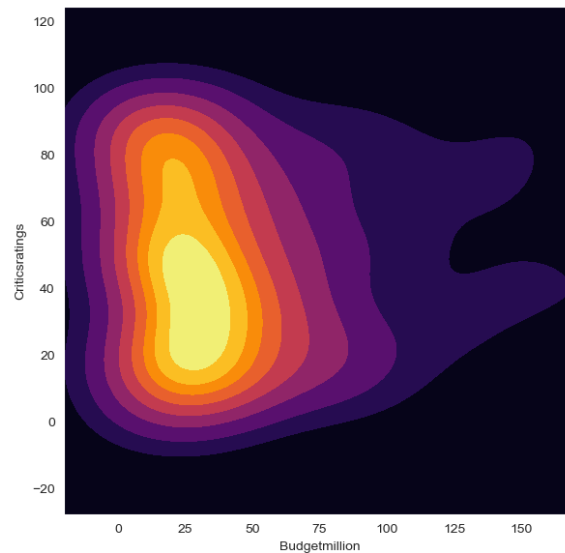
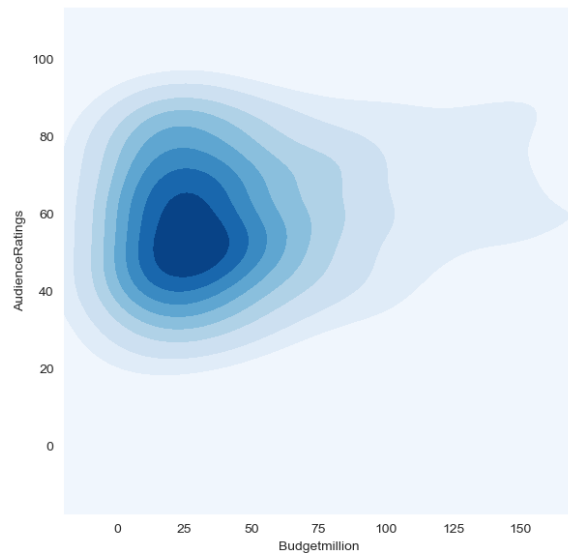


In [174]:

```

1  ##styling Dashboard
2  sns.set_style("darkgrid")
3  f,axes=plt.subplots(2,2,figsize=(15,15))
4  vis3 = sns.kdeplot(data=df, x="Budgetmillion", y="AudienceRatings",shade=True,sh
5  vis4 = sns.kdeplot(data=df, x="Budgetmillion", y="Criticsratings",shade=True,sha
6  vis3.set(xlim=(-20,170))
7  vis4.set(xlim=(-20,170))
8  vl=sns.violinplot(data=Thriller,x="Year",y="Criticsratings",ax=axes[1,0])
9  axes[1,1].hist(df[df.Genre=="Drama"].Budgetmillion)
10 plt.show()

```



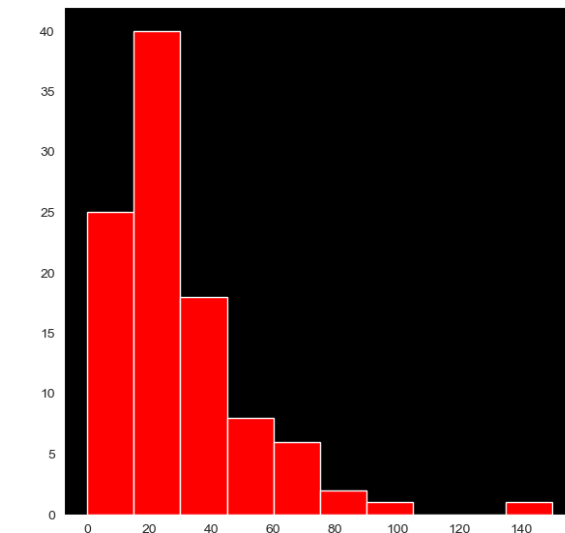
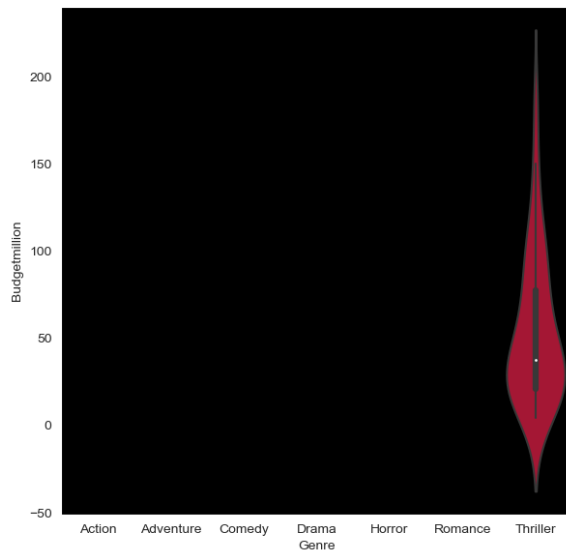
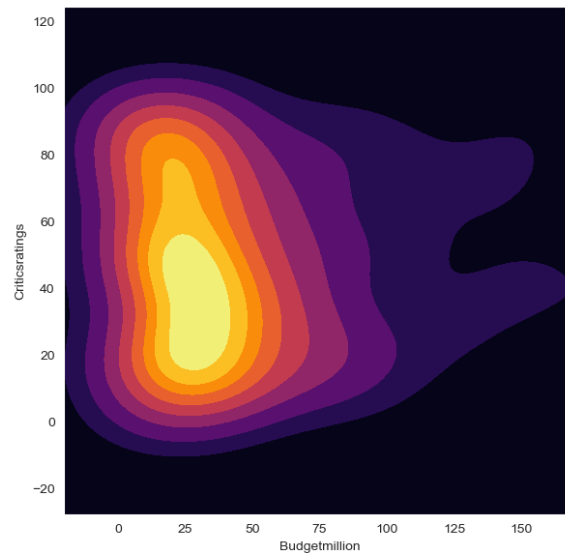
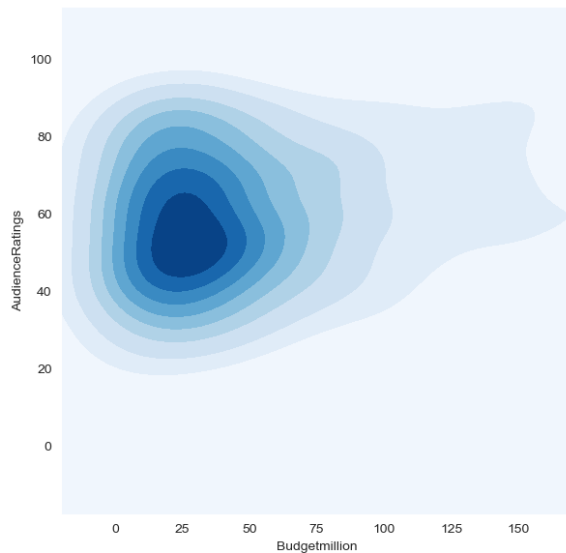
so in here , you can not change the background color of last two, if you do , you can pass the kewrod dictionaries to set style

In [187]:

```

1 sns.set_style("dark", {"axes.facecolor": "black"})
2 f, axes = plt.subplots(2, 2, figsize=(15, 15))
3 vis3 = sns.kdeplot(data=df, x="Budgetmillion", y="AudienceRatings", shade=True, s
4 vis4 = sns.kdeplot(data=df, x="Budgetmillion", y="Criticsratings", shade=True, s
5 vis3.set(xlim=(-20, 170))
6 vis4.set(xlim=(-20, 170))
7 vl = sns.violinplot(data=Thriller, x="Genre", y="Budgetmillion", ax=axes[1, 0], p
8 axes[1, 1].hist(df[df.Genre == "Drama"].Budgetmillion,color="red")
9 plt.show()
10

```

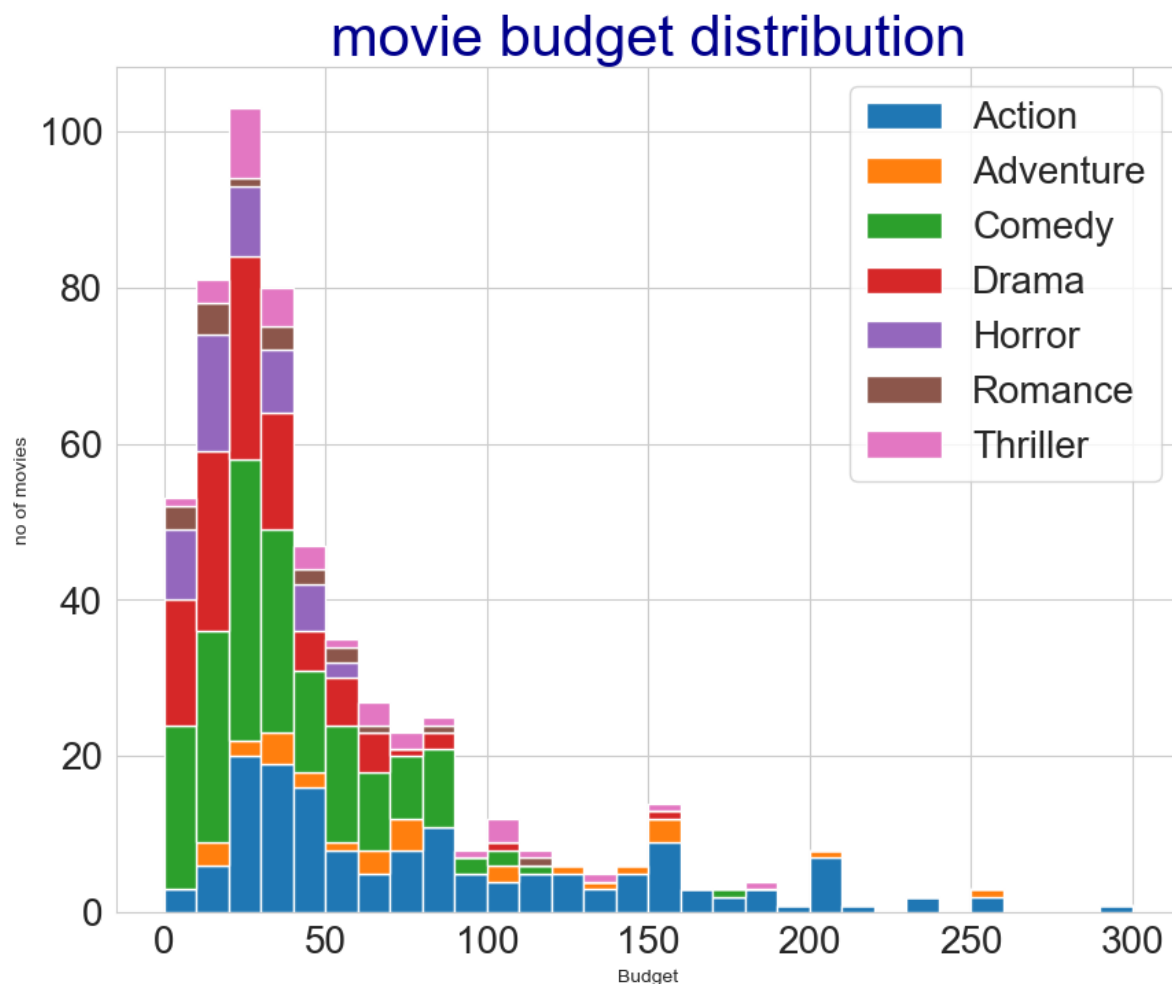


In [207]:

```

1  #improving the chart
2  import seaborn as sns
3  import matplotlib.pyplot as plt
4  list1 = list()
5  mylabels=list()
6
7  sns.set_style("whitegrid")
8  f, axes = plt.subplots()
9  f.set_size_inches(10, 8)
10 for gen in df.Genre.cat.categories:
11     list1.append(df[df.Genre == gen].Budgetmillion)
12     mylabels.append(gen)
13 plt.hist(list1,bins=30,stacked=True,rwidth=1,label=mylabels)
14 plt.title("movie budget distribution",fontsize=30,color="darkblue")
15 plt.ylabel("no of movies")
16 plt.xlabel("Budget")
17 plt.xticks(fontsize=20)
18 plt.yticks(fontsize=20)
19 plt.legend(fontsize=20)
20 plt.show()

```



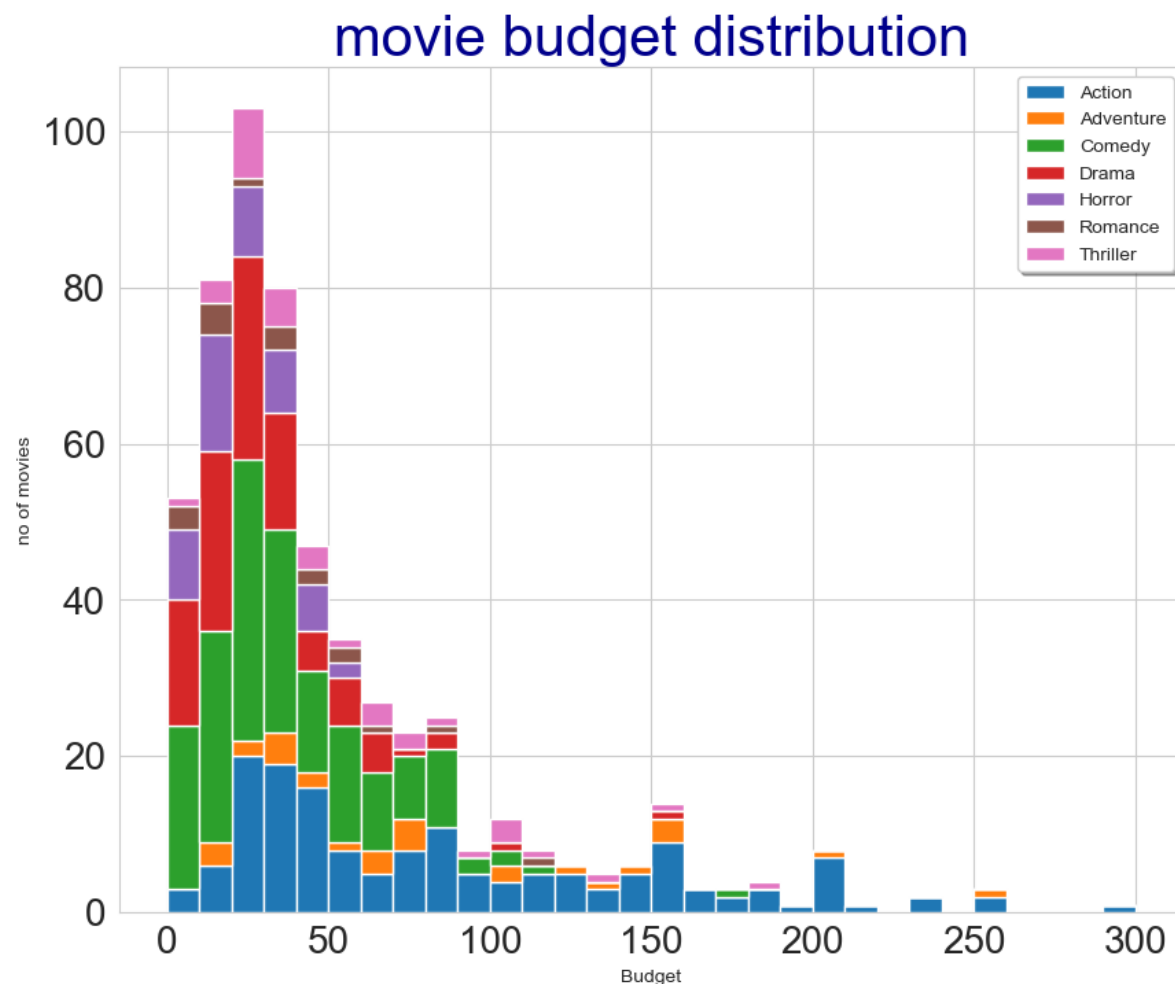


In [209]:

```

1 #improving the chart
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 list1 = list()
5 mylabels=list()
6
7 sns.set_style("whitegrid")
8 f, axes = plt.subplots()
9 f.set_size_inches(10, 8)
10 for gen in df.Genre.cat.categories:
11     list1.append(df[df.Genre == gen].Budgetmillion)
12     mylabels.append(gen)
13 plt.hist(list1,bins=30,stacked=True,rwidth=1,label=mylabels)
14 plt.title("movie budget distribution",fontsize=30,color="darkblue")
15 plt.ylabel("no of movies")
16 plt.xlabel("Budget")
17 plt.xticks(fontsize=20)
18 plt.yticks(fontsize=20)
19 plt.legend(shadow=True)
20 plt.show()

```



In [ ]:

1

