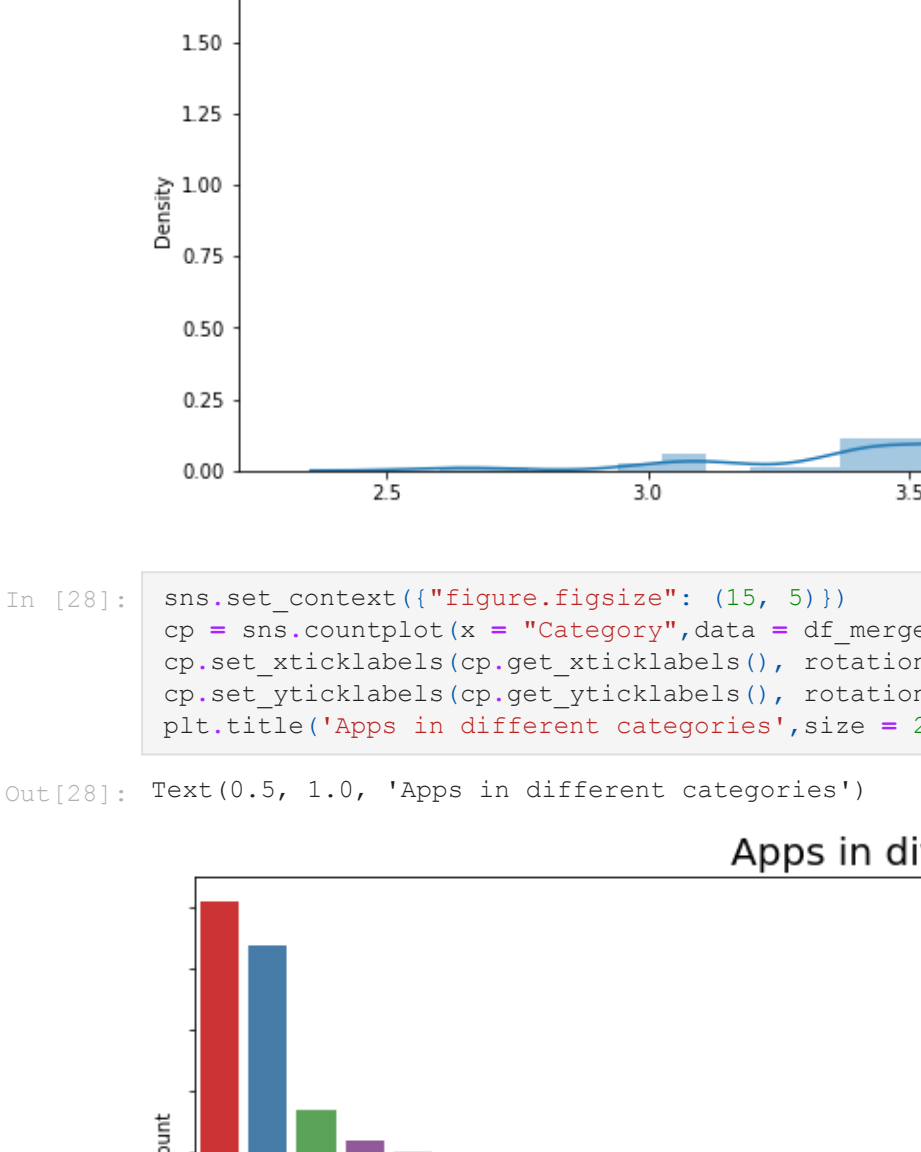
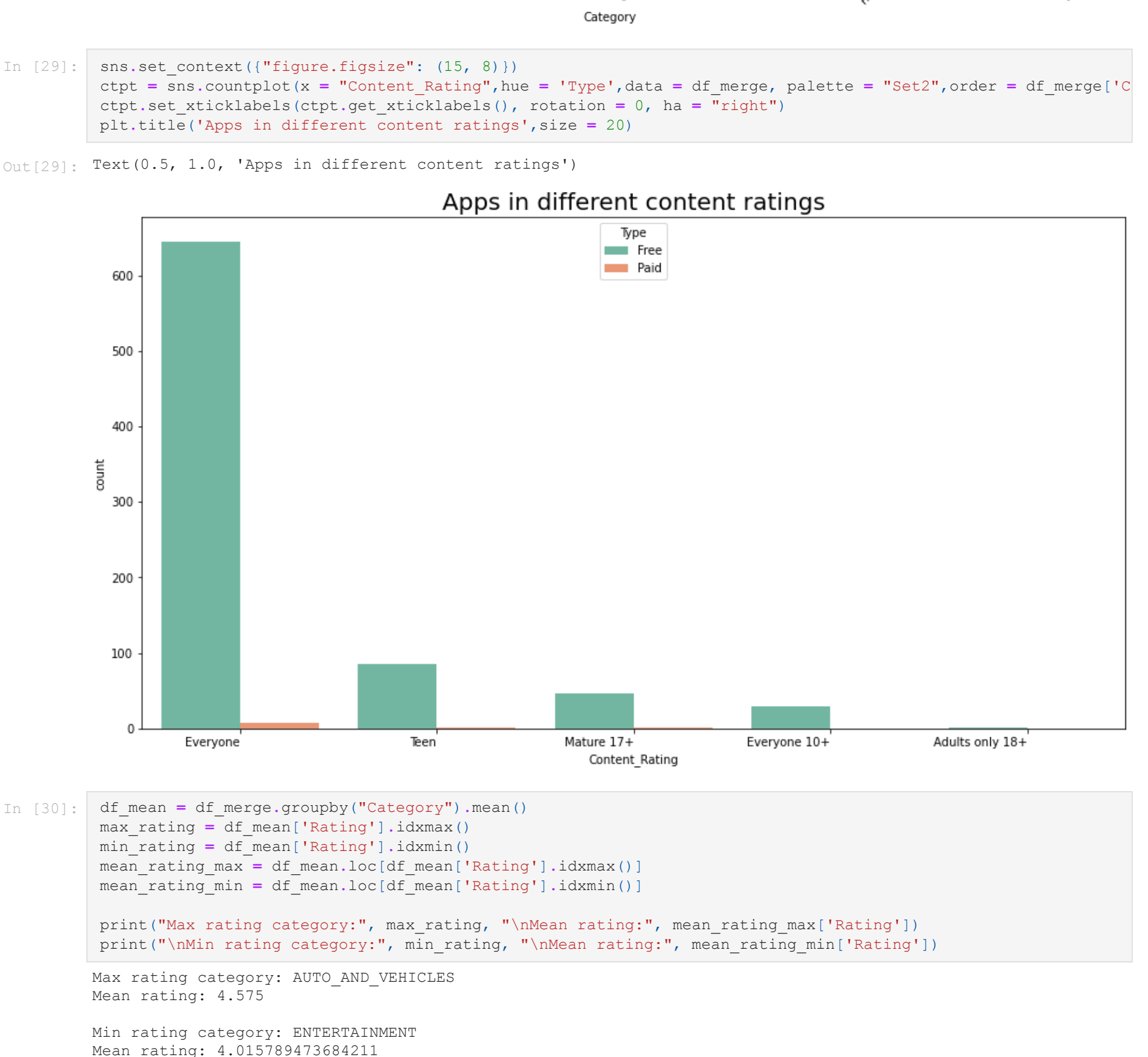


Common Trends

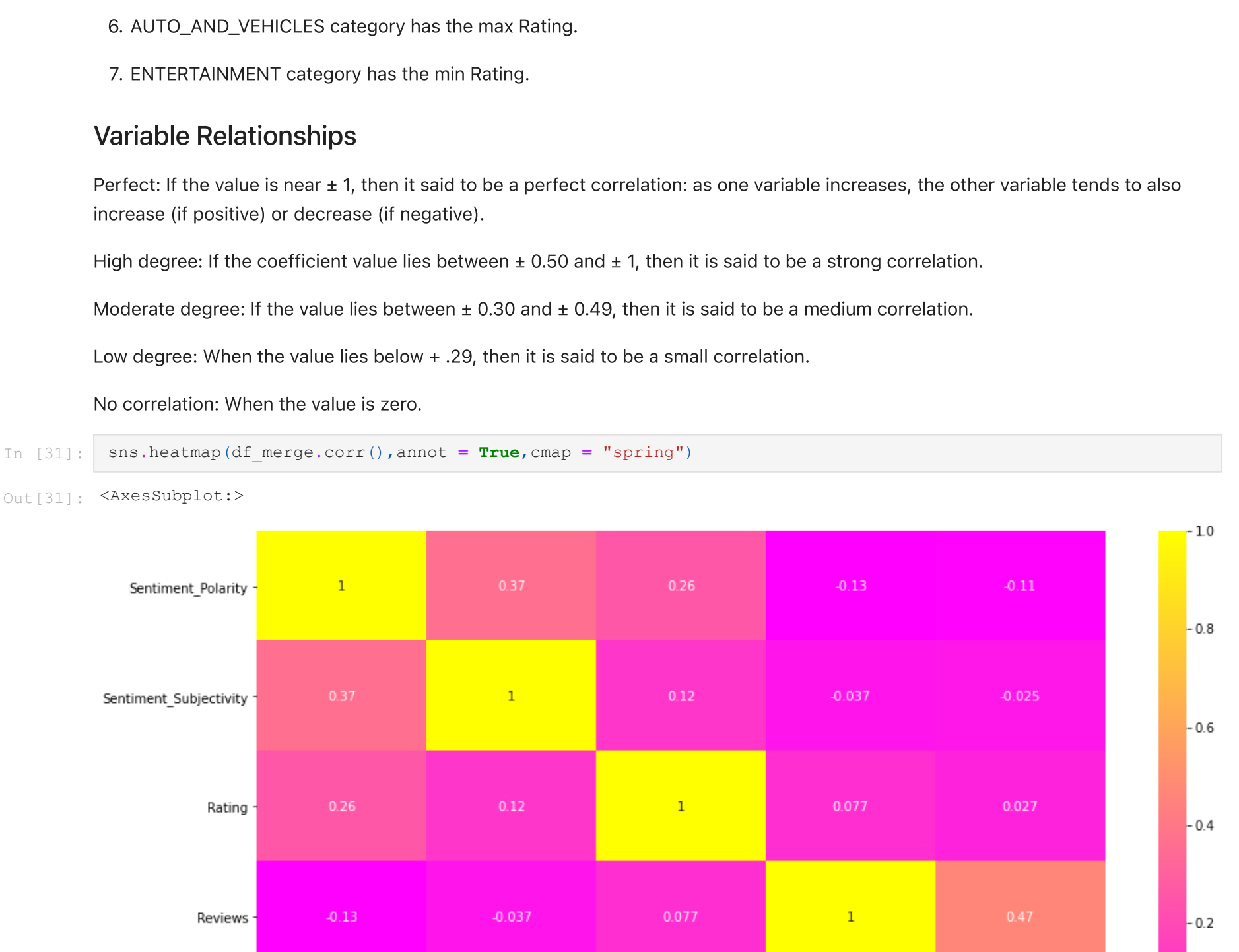
```
In [27]: plt.figure(figsize = (15,5))
sns.distplot(df_merge["Rating"])
```



```
In [28]: sns.set_context({"figure.figsize": (15, 5)})
cp = sns.countplot(x = "Category", data = df_merge, palette = "Set1", order = df_merge["Category"].value_counts().index)
cp.set_xticklabels(cp.get_xticklabels(), rotation = 45, ha = "right")
cp.set_yticklabels(cp.get_yticklabels(), rotation = 0, ha = "right")
plt.title('Apps in different categories',size = 20)
```



```
In [29]: sns.set_context({"figure.figsize": (15, 8)})
cp = sns.countplot(x = "Content_Rating", hue = "Type", data = df_merge, palette = "Set1", order = df_merge["Content_Rating"].value_counts().index)
cp.set_xticklabels(cp.get_xticklabels(), rotation = 0, ha = "right")
plt.title('Apps in different content ratings',size = 20)
```



```
In [30]: df_mean = df_merge.groupby("Category").mean()
max_rating = df_mean["Rating"].idxmax()
min_rating = df_mean["Rating"].idxmin()
mean_rating_max = df_mean.loc[df_mean["Rating"].idxmax()]
mean_rating_min = df_mean.loc[df_mean["Rating"].idxmin()]

print("Max rating category", max_rating, "\nMean rating", mean_rating_max["Rating"])
print("Min rating category", min_rating, "\nMean rating", mean_rating_min["Rating"])

Max rating category: AUTO_AND_VEHICLES
Min rating category: ENTERTAINMENT
Mean rating: 4.019789471082211
```

Conclusions:

- Most of the Ratings are in the range of 4.0 to 4.7.
- Most of the Apps belong to the category of Family.
- Category of Comics has the least Apps.
- Most of the Apps belongs to the content rating of Everyone.
- Most of the Apps belongs to the type of Free.
- AUTO_AND_VEHICLES category has the max Rating.
- ENTERTAINMENT category has the min Rating.

Variable Relationships

Perfect: If the value is near ± 1 , then it will be a perfect correlation: as one variable increases, the other variable tends to also increase (if positive) or decrease (if negative).

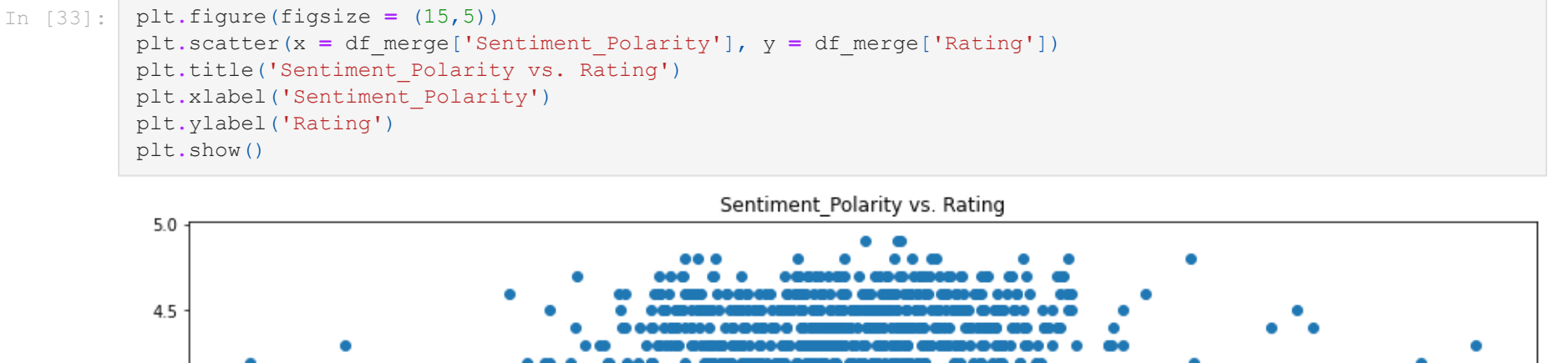
High degree: If the coefficient value lies between ± 0.50 and ± 1 , then it is said to be a strong correlation.

Moderate degree: If the value lies between ± 0.30 and ± 0.49 , then it is said to be a medium correlation.

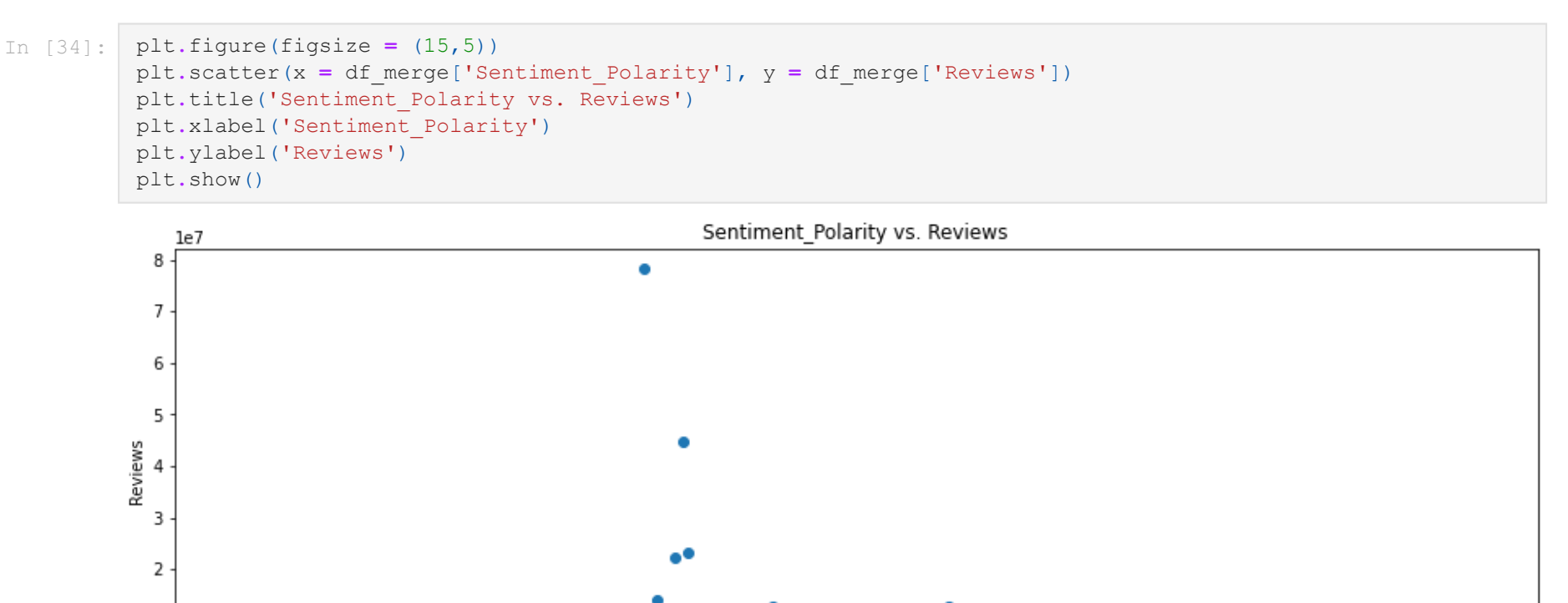
Low degree: When the value lies below $\pm .29$, then it is said to be a small correlation.

No correlation: When the value is zero.

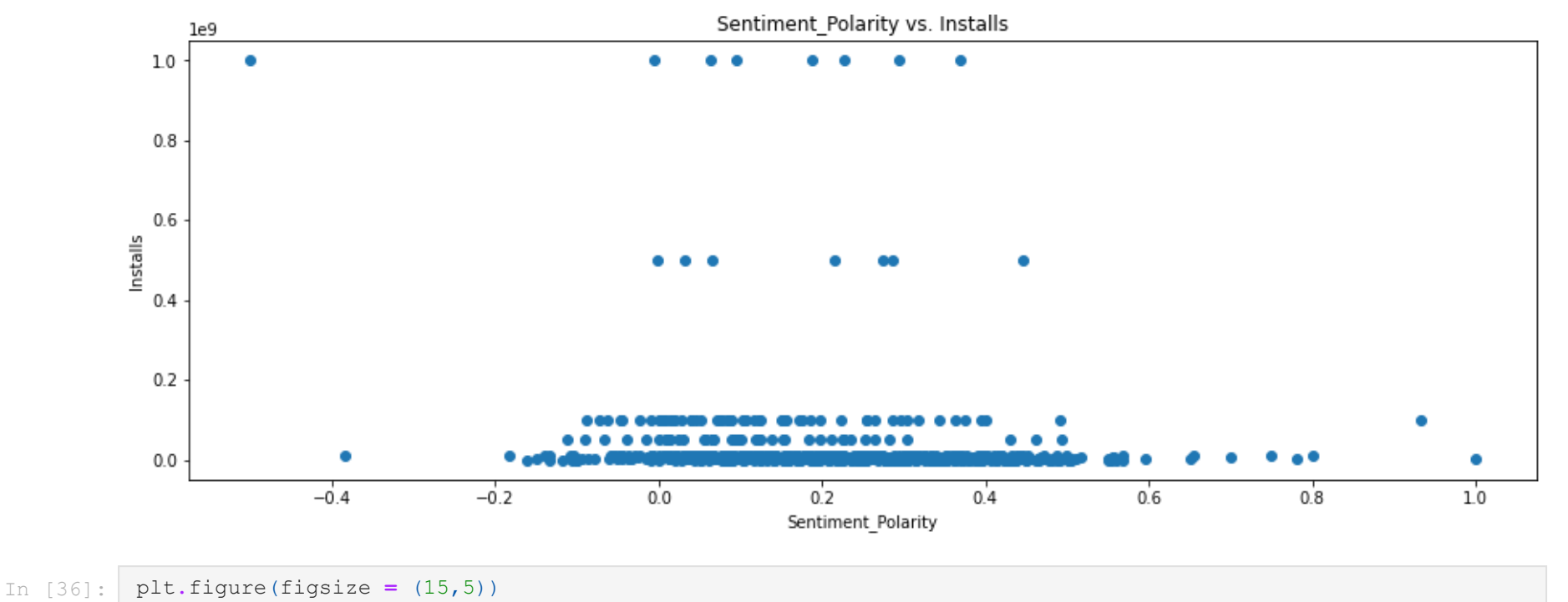
```
In [31]: sns.heatmap(df_merge.corr(),annot = True,cmap = "magma")
```



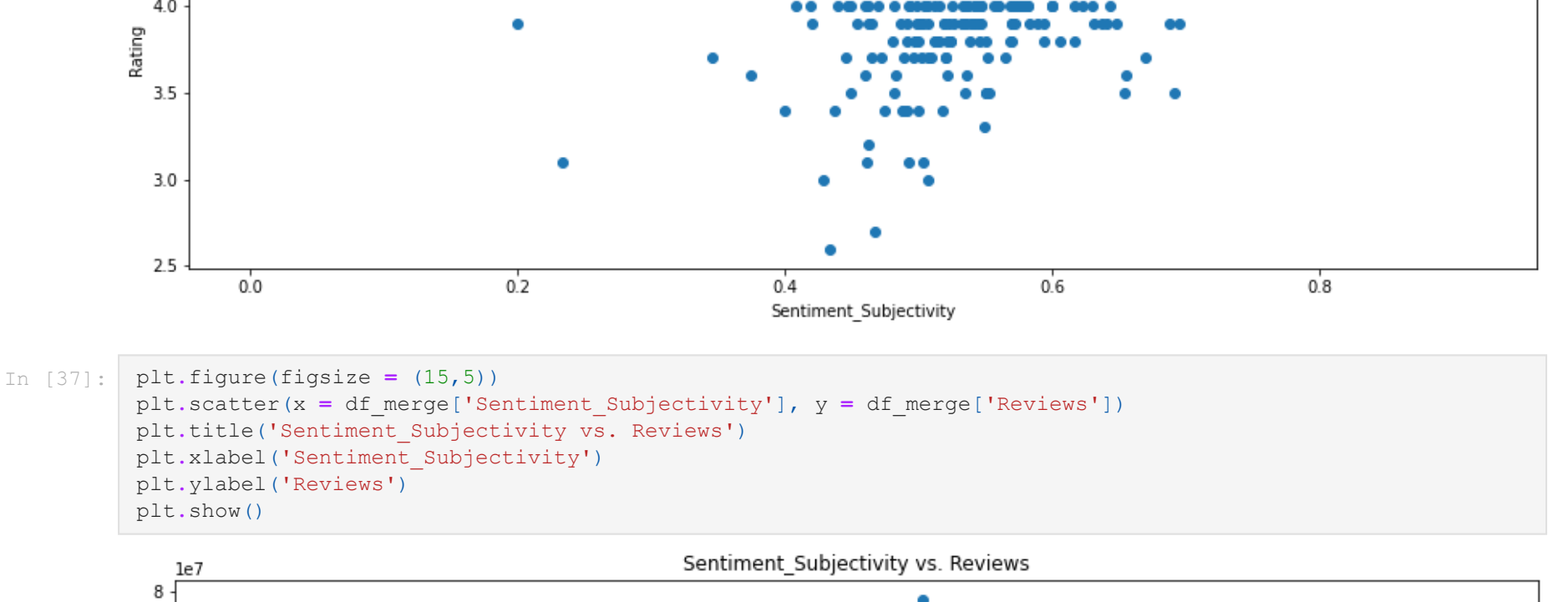
```
In [32]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Sentiment_Polarity"], y = df_merge["Sentiment_Subjectivity"])
plt.title('Sentiment_Polarity vs. Sentiment_Subjectivity')
plt.xlabel('Sentiment_Polarity')
plt.ylabel('Sentiment_Subjectivity')
plt.show()
```



```
In [33]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Sentiment_Polarity"], y = df_merge["Rating"])
plt.title('Sentiment_Polarity vs. Rating')
plt.xlabel('Sentiment_Polarity')
plt.ylabel('Rating')
plt.show()
```



```
In [34]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Sentiment_Polarity"], y = df_merge["Reviews"])
plt.title('Sentiment_Polarity vs. Reviews')
plt.xlabel('Sentiment_Polarity')
plt.ylabel('Reviews')
plt.show()
```



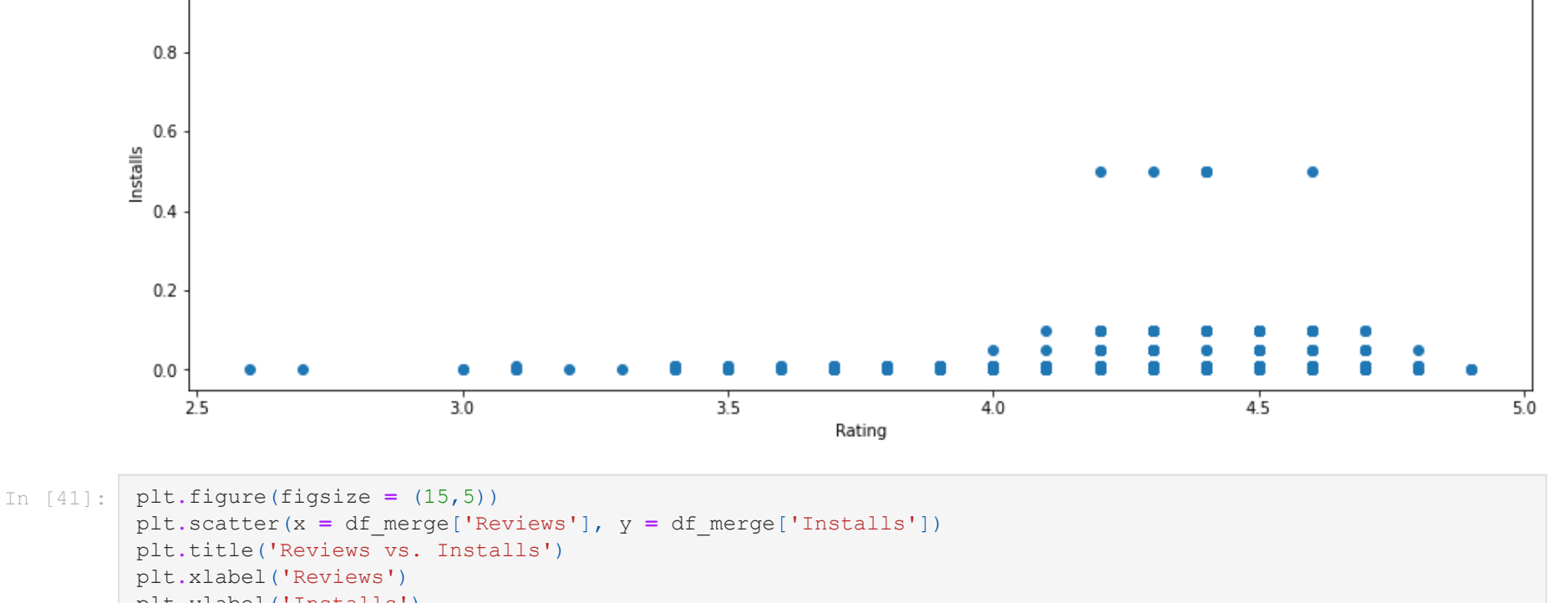
```
In [35]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Sentiment_Polarity"], y = df_merge["Installs"])
plt.title('Sentiment_Polarity vs. Installs')
plt.xlabel('Sentiment_Polarity')
plt.ylabel('Installs')
plt.show()
```



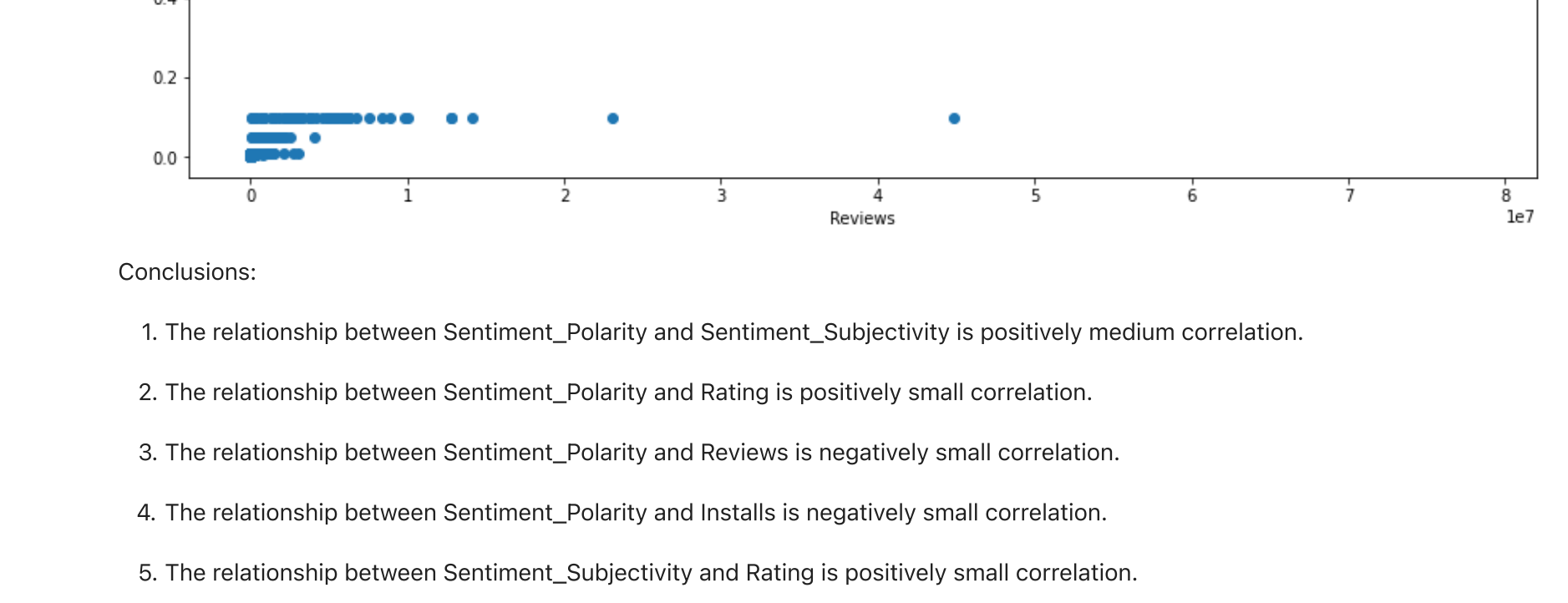
```
In [36]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Sentiment_Subjectivity"], y = df_merge["Rating"])
plt.title('Sentiment_Subjectivity vs. Rating')
plt.xlabel('Sentiment_Subjectivity')
plt.ylabel('Rating')
plt.show()
```



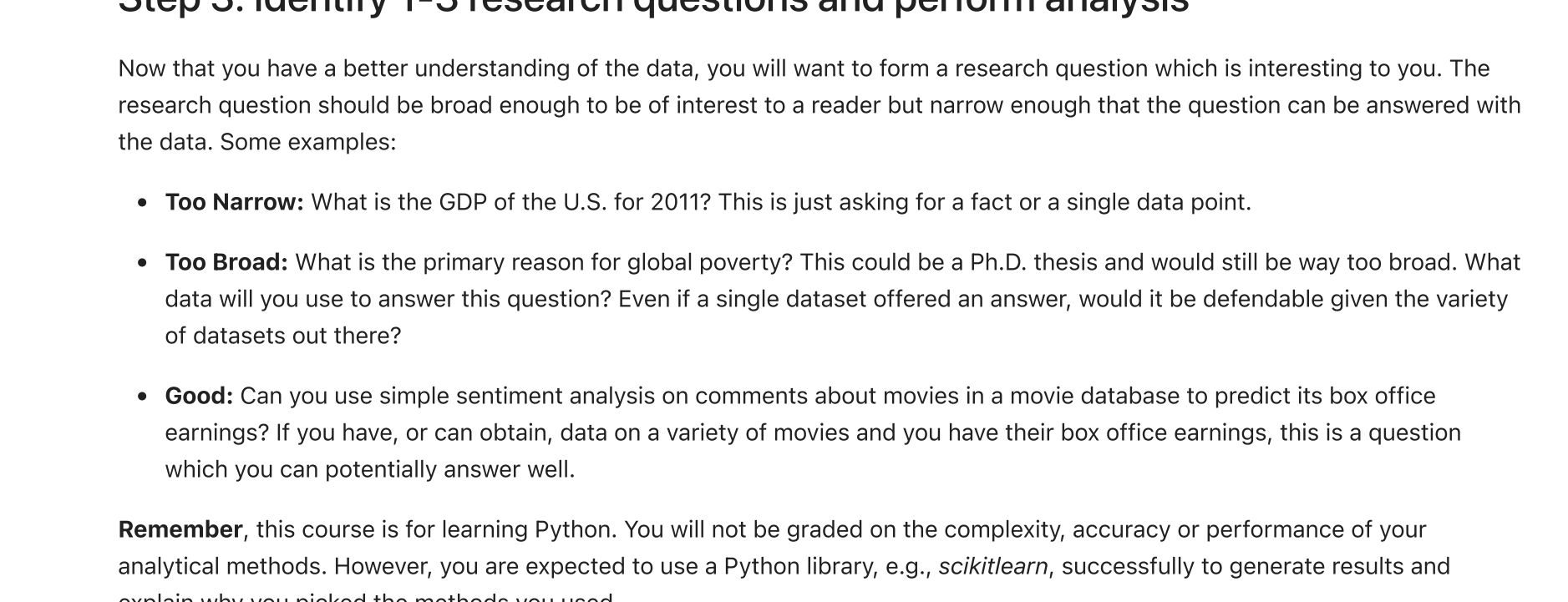
```
In [37]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Sentiment_Subjectivity"], y = df_merge["Reviews"])
plt.title('Sentiment_Subjectivity vs. Reviews')
plt.xlabel('Sentiment_Subjectivity')
plt.ylabel('Reviews')
plt.show()
```



```
In [38]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Sentiment_Subjectivity"], y = df_merge["Installs"])
plt.title('Sentiment_Subjectivity vs. Installs')
plt.xlabel('Sentiment_Subjectivity')
plt.ylabel('Installs')
plt.show()
```



```
In [39]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Rating"], y = df_merge["Installs"])
plt.title('Rating vs. Installs')
plt.xlabel('Rating')
plt.ylabel('Installs')
plt.show()
```



```
In [40]: plt.figure(figsize = (15,5))
plt.scatter(x = df_merge["Reviews"], y = df_merge["Installs"])
plt.title('Reviews vs. Installs')
plt.xlabel('Reviews')
plt.ylabel('Installs')
plt.show()
```



Conclusions:

- The relationship between Sentiment_Polarity and Sentiment_Subjectivity is positively medium correlation.
- The relationship between Sentiment_Polarity and Rating is positively small correlation.
- The relationship between Sentiment_Polarity and Reviews is negatively small correlation.
- The relationship between Sentiment_Polarity and Installs is negatively small correlation.
- The relationship between Sentiment_Subjectivity and Rating is positively small correlation.
- The relationship between Sentiment_Subjectivity and Reviews is negatively small correlation.
- The relationship between Sentiment_Subjectivity and Installs is negatively small correlation.
- The relationship between Rating and Reviews is positively small correlation.
- The relationship between Rating and Installs is positively small correlation.
- The relationship between Reviews and Installs is positively medium correlation.

Step 3: Identify 1-3 research questions and perform analysis

Now that you have a better understanding of the data, you will want to form a research question which is interesting to you. The research question should be broad enough to be of interest to a reader but narrow enough that the question can be answered with the data. Some examples:

- Too Narrow:** What is the GDP of the U.S. for 2012? This is just asking for a fact or a single data point.
- Too Broad:** What is the primary reason for global poverty? This could be a Ph.D. thesis and would be way too broad. What data will you use to answer this question? Even if a single dataset offered an answer, it would be defendable given the variety of datasets out there.
- Good:** Can you use simple sentiment analysis on comments about movies in a movie database to predict its box office earnings? If you have, or can obtain, data on a variety of movies and you have their box office earnings, this is a question which you can potentially answer well.

Remember, this course is for learning Python. You will not be graded on the complexity, accuracy or performance of your analytical methods. However, you are expected to use a Python library, e.g., `scikit-learn`, successfully to generate results and explain why you picked the methods you used.

Research Questions

- Can we predict App Ratings for Apps on the Google Play Store?
- What categories of Apps that a marketing strategy should targeted on the Google Play Store Apps?
- Can we use NLP to train the Google Play Store dataset to predict reviews?

Perform Analysis

Question 1

1. Can we predict App Ratings for Apps on the Google Play Store?

```
In [42]: data = df_merge.copy()
data = data.drop(["App_name", "Category", "Size", "Type", "Price", "Genres", "Content_Rating", "Last_Updated"])

In [43]: features = ["Sentiment_Polarity", "Sentiment_Subjectivity", "Reviews", "Installs"]
X = data[features]
y = data["Rating"]
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20, random_state = 10)
```

Linear Regression

```
In [44]: # Linear Regression: Fit a model to the training set
regressor = LinearRegression()
regressor.fit(X_train, y_train)

# Perform Prediction using Linear Regression Model
y_prediction = regressor.predict(X_test)
```

```
Out[44]: array([4.25553971, 4.2286547 , 4.35424716, 4.42727049, 4.2310485 ,
4.47165403, 4.20162923, 4.22941127, 4.26399847, 4.36486569,
4.29835892, 4.26043216, 4.25357674, 4.26607465, 4.2189732 ,
4.3501017 , 4.26569116, 4.21022499, 4.21008285, 4.35021077,
4.29058251, 4.29007471, 4.4226582 , 4.27598018, 4.27196247,
4.31872617, 4.20829212, 4.16648117, 4.33784355, 4.18151303,
4.30394538, 4.25234116, 4.24943659, 4.2597446 , 4.2205985 ,
4.27932197, 4.28048079, 4.28537028, 4.2624821 , 4.30671991,
4.3614113 , 4.22125539, 4.08118405, 4.2320141 , 4.3566559 ,
4.25474776, 4.2344962 , 4.34317373, 4.3551423 , 4.27196247,
4.14070219, 4.18976119, 4.1677009 , 4.1984759 , 4.26647002,
4.40513432, 4.28661551, 4.31598994, 4.21820463, 4.1202365 ,
4.45089407, 4.2653815 , 4.20791059, 4.26458727, 4.19584345,
4.33502164, 4.34948202, 4.29867393, 4.17535956, 4.3316468 ,
4.26144271, 4.16086531, 4.23012029, 4.23562605, 4.28548397,
4.404492 , 4.30374748, 4.20693939, 4.36145847, 4.35263155,
4.13303933, 4.40777004, 4.14487009, 4.17238862, 4.42831678,
4.1695992 , 4.22040364, 4.16214687, 4.28540037, 4.34818193 ,
4.23733801, 4.12157931, 4.10567099, 4.26484595, 4.17506204,
4.22134492, 4.37363853, 4.15302733, 4.4124037 , 4.16837848 ,
4.3289064 , 4.2602448 , 4.37534786, 4.39864578, 4.25981085,
4.25978715, 4.3325213 , 4.15998625, 4.24763972, 4.22592215,
4.30673529, 4.35962063, 4.27135825, 4.18066862, 4.26895179,
4.2664462 , 4.14774059, 4.34883314, 4.1162979 , 4.37294756,
4.210739 , 4.34811082, 4.31610055, 4.1139386 , 4.17893862,
4.33634793, 4.23946261, 4.22564516, 4.2073505 , 4.22216767,
4.3984686 , 4.14621244, 4.21516356, 4.26484595, 4.17506204,
4.21713998, 4.3403376 , 4.28675259, 4.26743107, 4.20464702,
4.27596925, 4.30319311, 4.09260822, 4.3797093 , 4.23623904,
4.2287795 , 4.19187434, 4.40202139, 4.31741042, 4.21975655,
4.07195389, 4.47481739, 4.38179984, 4.18871833, 4.25762449,
4.38054068, 4.33685774, 4.20717396, 4.28672162, 4.35286683,
4.43336102, 4.26845055, 4.17113241, 4.18391051])
```

```
In [45]: # The mean of the expected target value in test set
y_test.describe()
```

```
Out[45]: count    164.000000
mean      4.259756
std       0.328544
min       2.600000
25%      4.100000
50%      4.300000
75%      4.500000
max       4.900000
Name: Rating, dtype: float64
```

```
In [46]: # Evaluate Linear Regression Accuracy using Root Mean Square Error
RMSE_Linear_Regression = sqrt(mean_squared_error(y_true = y_test, y_pred = y_prediction))
print(RMSE_Linear_Regression)

0.319534285398724
```

Decision Tree Regression

```
In [47]: # Decision Tree Regression: Fit a new regression model to the training set
regressor = DecisionTreeRegressor(max_depth=10)
regressor.fit(X_train, y_train)

# Perform Prediction using Decision Tree Regressor
y_prediction = regressor.predict(X_test)
```

```
Out[47]: array([3.5 , 4.6 , 4.1 , 4.53333333, 4.35 ,
4.38 , 4.6 , 4.39375 , 4.5 , 4.35263158,
4.41428571, 4.6 , 4.5263158 , 4.3 , 4.5 ,
4.1 , 4.3 , 4.3 , 4.2 , 4.2608696 , 4.17893862,
4.33333333, 4.29 , 4.53333333, 4.33333333, 3.9 ,
4.22658251, 4.21428571, 4.21428571, 4.263158 , 4.6 ,
4.6 , 4.3 , 4.5 , 3.7 , 3.5 ,
4.3 , 4.33333333, 4. , 4. ,
4.6 , 3.8 , 4.2608696 , 4.2608696 ,
4.9 , 4.3 , 4.22608696 , 4.56666667,
3.6 , 4.3 , 4.5263158 , 3.825 , 4.33333333,
4.1 , 4.41428571, 4.15333333, 4.3 ,
4.53333333, 4.275 , 3.9 , 4.1 , 3.99090909,
4.22608696, 4.6 , 4.15 , 3.8 ,
4.3 , 3.99090909, 4.53333333, 4.4 , 4.5 ,
4.6 , 4.6 , 4.175 , 4.35263158, 4.5 ,
4.3 , 4.3 , 4.3 , 4.2 , 4.2608696 ,
4.3 , 4.6 , 4.3 , 4.2 , 4.35263158,
3.875 , 3.1 , 4.38 , 3.5 , 4.53333333,
4.16 , 4.08 , 4.3 , 4.4 , 4.5 ,
4.2 , 4.35 , 4.08 , 4.3 , 4.5 ,
4.5 , 4.15 , 4.2 , 3.9 ,
4.33333333, 4.08 , 4.175 , 3.9 ,
4.3 , 4.6 , 4.4 , 4.4 , 4.36666667,
4.7 , 4.2 , 4.41428571, 4.21428571, 4.1 ,
4.15 , 4.2 , 4.21428571, 4.21428571, 4.33333333,
4.58 , 4.17777778, 4.3 , 4.53333333, 4.15 ,
4.1 , 4.6 , 4.425 , 4.3 , 4.56666667,
4.175 , 3.5 , 3. , 4.08 , 4.1 ,
4.1 , 4.6 , 4.66666667, 3.8 , 4.1 ,
3.4 , 4.4 , 4.4 , 3.7 , 4.4 ,
4.53333333, 4.33333333, 4.6 , 4.36666667, 4.2 ,
4.3 , 4.275 , 3.9090909 , 4.1 )
```

```
In [48]: # The mean of the expected target value in test set
y_test.describe()
```

```
Out[48]: count    164.000000
mean      4.259756
std       0.328544
min       2.600000
25%      4.100000
50%      4.300000
75%      4.500000
max       4.900000
Name: Rating, dtype: float64
```

```
In [49]: # Evaluate Decision Tree Regression Accuracy using Root Mean Square Error
RMSE_Decision_Tree_Regression = sqrt(mean_squared_error(y_true = y_test, y_pred = y_prediction))
print(RMSE_Decision_Tree_Regression)

0.292674208912836
```

Random Forest Regression

```
In [50]: # Random Forest Regression: Fit a new regression model to the training set
model = RandomForestRegressor(n_estimators = 200, n_jobs = -1, random_state = 10)
model.fit(X_train, y_train)

# Perform Prediction using Random Forest Regressor
y_prediction = model.predict(X_test)
```

```
Out[50]: array([4.254 , 4.518 , 4.43 , 4.427 , 4.371 , 4.362 , 4.507 , 4.383 ,
4.4005, 4.3075, 4.383 , 4.463 , 4.276 , 4.4635, 4.4753, 4.604 ,
4.2665, 4.587 , 3.9445, 4.5375, 4.32 , 4.1085, 4.4535, 4.407 ,
4.179 , 4.354 , 4.019 , 4.1035, 4.279 , 4.132 , 4.489 , 4.184 ,
4.4285, 3.956 , 4.1365, 4.5405, 4.261 , 4.3845, 4.1815, 4.4955,
4.366 , 4.508 , 4.0065, 4.3755, 4.3095, 4.4245, 4.4565, 4.0625,
4.1825, 4.245 , 4.4425, 4.0365, 4.4385, 4.2395, 3.457 , 4.149 ,
4.128 , 4.2825, 4.4535, 4.27 , 4.026 , 4.2255, 4.321 , 4.0365,
4.486 , 4.2745, 4.3145, 4.523 , 4.504 , 4.198 , 4.3225, 4.33 ,
4.231 , 4.3285, 4.0345, 4.188 , 4.3335, 3.6865, 4.452 , 4.408 ,
4.214 , 4.2805, 4.1605, 3.828 , 4.395 , 4.244 , 4.5625, 4.008 ,
4.2295, 4.345 , 4.444 , 4.378 , 4.3665, 4.4693, 4.297 , 4.415 ,
4.295 , 4.465 , 4.1155, 4.1555, 4.264595 , 4.17506204,
4.21713998, 4.3403376 , 4.28675259, 4.26743107, 4.20464702,
4.27596925, 4.30319311, 4.09260822, 4.3797093 , 4.23623904,
4.2287795 , 4.19187434, 4.40202139, 4.31741042, 4.21975655,
4.07195389, 4.47481739, 4.38179984, 4.18871833, 4.25762449,
4.38054068, 4.33685774, 4.20717396, 4.28672162, 4.35286683,
4.43336102, 4.26845055, 4.17113241, 4.18391051])
```

```
In [51]: # The mean of the expected target value in test set
y_test.describe()
```

```
Out[51]: count    164.000000
mean      4.259756
std       0.328544
min       2.600000
25%      4.100000
50%      4.300000
75%      4.500000
max       4.900000
Name: Rating, dtype: float64
```

```
In [52]: # Evaluate Random Forest Regression Accuracy using Root Mean Square Error
RMSE_Random_Forest = sqrt(mean_squared_error(y_true = y_test, y_pred = y_prediction))
print(RMSE_Random_Forest)

0.292674208912836
```

Question 2

1. What categories of Apps that a marketing strategy should targeted on the Google Play Store Apps?

```
In [53]: data_clustering = df_merge.copy()
data_clustering = data_clustering.drop(["App_name", "Category", "Size", "Type", "Price", "Genres", "Content_Rating", "Last_Updated"])
data_clustering = preprocessing.normalize(data_clustering)
print("The normalized data for clustering is: ", data_clustering)
```

```
X_input = StandardScaler().fit_transform(normalized_data)

# plot the elbow method
plt.figure(figsize = (15,8))
plt.subplot(2, 1, 1):
for k in range(1, 11):
    kmeans = KMeans(n_clusters = k, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
    wcss = kmeans.inertia_
    plt.plot(range(1, 11), wcss)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WCSS')
plt.grid()
plt.tight_layout()
plt.show()
```

```
# find the best elbow
k1 = KMeans(range(1, 11), wcss, curve = "convex", direction = "decreasing")
elbow = k1.elbow
print("The best elbow (n_clusters) is", elbow, "\n")

# print the model, fit the model to the normalized data
kmeans = KMeans(n_clusters = elbow)
model = kmeans.fit(X_input)
print("Model: k1, model: \n")

# print the resulting cluster centers
centers = model.cluster_centers_
print("Cluster centers: \n", centers, "\n")

predict = kmeans.predict(X_input)
```

```
plt.figure(figsize = (15,8))
plt.subplot(2, 1, 2)
plt.scatter(X_input[:, 0], X_input[:, 1], c = predict, s = 50, cmap = "viridis")
plt.scatter(centers[:, 0], centers[:, 1], c = "black", s = 200, alpha = 0.5)
plt.grid()
```

```
The normalized data for clustering:
[[19.39397124e-07 1.11902073e-06 7.99990080e-06 4.97993825e-03
1.56995674e-08 5.12348310e-08 3.79994548e-07 4.87314214e-02
1.49988166e-07 6.034849e-07 2.62307192e-06 1.03224500e-02
6.99944722e-01
1.93524188e-10 5.09290960e-09 4.29487722e-08 4.87982515e-02
5.98808655e-01
1.98864210e-07 5.46797850e-07 4.29377090e-06 1.03224500e-02
5.99944722e-01
1.71851913e-09 4.30336705e-07 4.09835950e-06 2.28269730e-02
4.93987613e-01]]
```


The best elbow (n_clusters) is 3

```
Model:
KMeans(n_clusters=3)
```

```
Cluster centers:
[[ -0.03587705 -0.03454877 -0.03492727 -0.06608097 0.08384201]
[ 2.61622057  6.12362753  26.23071922  6.79976019 -13.1148128]
[ 0.33320495  0.25569829  0.28401794  4.45455522  7.80421911]]
```



```
In [54]: # Function that creates a DataFrame with a column for Cluster Number

def pd_centers(featuresUsed, centers):
    colNames = list(featuresUsed)
    colNames.append("Cluster")

    # Zip with a column called 'cluster' (index)
    Z = [np.append(i, index) for index, i in enumerate(centers)]

    # Convert to pandas data frame for plotting
    P = pd.DataFrame(Z, columns = colNames)
    P["Cluster"] = P["Cluster"].astype(int)
    return P
```

```
In [55]: features = ["Sentiment_Polarity", "Sentiment_Subjectivity", "Reviews", "Installs", "Rating"]
P = pd_centers(features, centers)
print(P)
```

```
0 Sentiment_Polarity Sentiment_Subjectivity Reviews Installs Rating \
1 0.03980777 0.203449 0.034927 -0.06608097 0.083842
2 2.616221 6.123628 26.233109 6.799761 -13.114813
3 0.333205 0.255968 0.284018 4.454555 -7.804219
```

```
0 Cluster
1 1
2 2
```

```
In [56]: # Function that creates Parallel Plots

def parallel_plot(data,
my_colors = list(islice(cycle(['b', 'r', 'g', 'y', 'k']), None, len(data))),
plt.figure(figsize = (15,8)).gca().axes.set_ylim([-20,30])
parallel_coordinates(data, 'Cluster', color = my_colors, marker = 'o')
```

```
In [57]: parallel_plot(P)
```



```
In [58]: new_data = data_clustering.copy()
predict_data = kmeans.predict(X_input)
new_data["Cluster"] = predict_data
new_data["Category"] = df_merge.Category

new_data.shape = new_data.shape
print("DataFrame's shape: \n", new_data.shape, "\n")

new_data.head = new_data.head()
print("DataFrame's Head: \n", new_data.head, "\n")
```